

22-冷链周转：HTTP的缓存代理

在[第20讲](#)中，我介绍了HTTP的缓存控制，[第21讲](#)我介绍了HTTP的代理服务。那么，把这两者结合起来就是这节课所要说的“缓存代理”，也就是支持缓存控制的代理服务。

之前谈到缓存时，主要讲了客户端（浏览器）上的缓存控制，它能够减少响应时间、节约带宽，提升客户端的用户体验。

但HTTP传输链路上，不只是客户端有缓存，服务器上的缓存也是非常有价值的，可以让请求不必走完整个后续处理流程，“就近”获得响应结果。

特别是对于那些“读多写少”的数据，例如突发热点新闻、爆款商品的详情页，一秒钟内可能有成千上万次的请求。即使仅仅缓存数秒钟，也能够把巨大的访问流量挡在外面，让RPS（request per second）降低好几个数量级，减轻应用服务器的并发压力，对性能的改善是非常显著的。

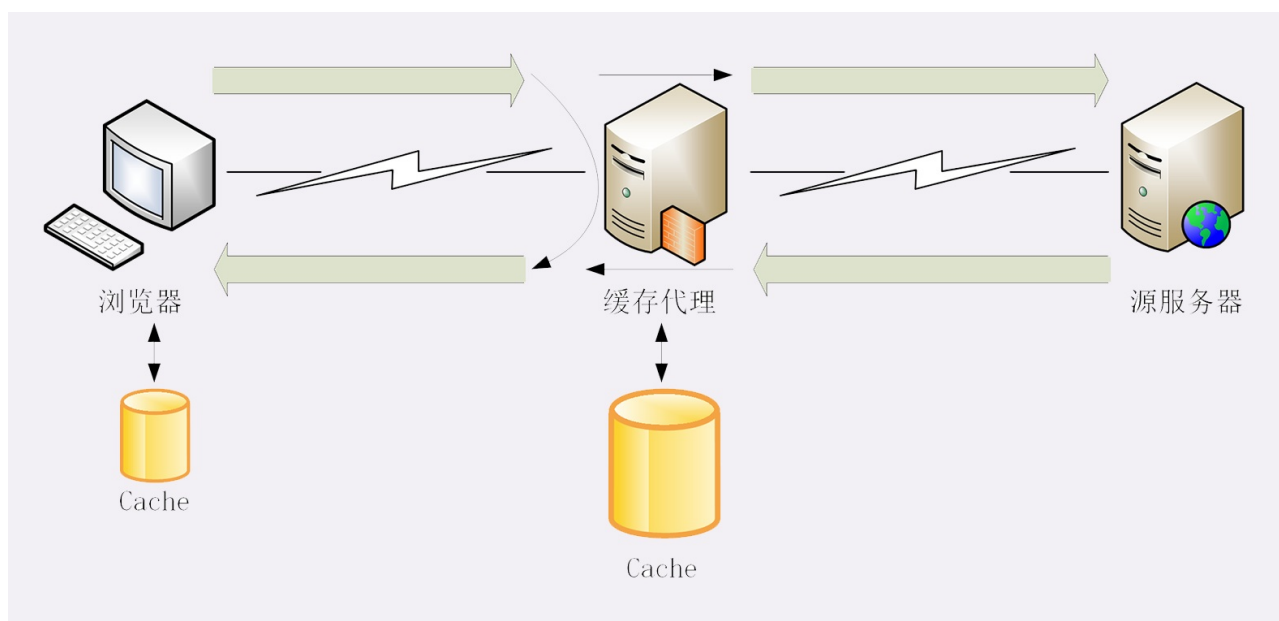
HTTP的服务器缓存功能主要由代理服务器来实现（即缓存代理），而源服务器系统内部虽然也经常有各种缓存（如Memcache、Redis、Varnish等），但与HTTP没有太多关系，所以这里暂且不说。

缓存代理服务

我还是沿用“生鲜速递+便利店”的比喻，看看缓存代理是怎么回事。

便利店作为超市的代理，生意非常红火，顾客和超市双方都对现状非常满意。但时间一长，超市发现还有进一步提升的空间，因为每次便利店接到顾客请求后都要专车跑一趟超市，还是挺麻烦的。

干脆这样吧，给便利店配发一个大冰柜。水果海鲜什么的都可以放在冰柜里，只要产品在保鲜期内，就允许顾客直接从冰柜提货。这样便利店就可以一次进货多次出货，省去了超市之间的运输成本。



通过这个比喻，你可以看到：在没有缓存的时候，代理服务器每次都是直接转发客户端和服务器的报文，中间不会存储任何数据，只有最简单的中转功能。

加入了缓存后就不一样了。

代理服务收到源服务器发来的响应数据后需要做两件事。第一个当然是把报文转发给客户端，而第二个就是把报文存入自己的Cache里。

下一次再有相同的请求，代理服务器就可以直接发送304或者缓存数据，不必再从源服务器那里获取。这样就降低了客户端的等待时间，同时节约了源服务器的网络带宽。

在HTTP的缓存体系中，缓存代理的身份十分特殊，它“既是客户端，又是服务器”，同时也“既不是客户端，又不是服务器”。

说它“即是客户端又是服务器”，是因为它面向源服务器时是客户端，在面向客户端时又是服务器，所以它即可以用客户端的缓存控制策略也可以用服务器端的缓存控制策略，也就是说它可以同时使用第20讲的各种“Cache-Control”属性。

但缓存代理也“即不是客户端又不是服务器”，因为它只是一个数据的“中转站”，并不是真正的数据消费者和生产者，所以还需要有一些新的“Cache-Control”属性来对它做特别的约束。

源服务器的缓存控制

[第20讲](#)介绍了4种服务器端的“Cache-Control”属性：max-age、no_store、no_cache和must-revalidate，你应该还有印象吧？

这4种缓存属性可以约束客户端，也可以约束代理。

但客户端和代理是不一样的，客户端的缓存只是用户自己使用，而代理的缓存可能会为非常多的客户端提供服务。所以，需要对它的缓存再多一些限制条件。

首先，我们要区分客户端上的缓存和代理上的缓存，可以使用两个新属性“private”和“public”。

“private”表示缓存只能在客户端保存，是用户“私有”的，不能放在代理上与别人共享。而“public”的意思就是缓存完全开放，谁都可以存，谁都可以用。

比如你登录论坛，返回的响应报文里用“Set-Cookie”添加了论坛ID，这就属于私人数据，不能存在代理上。不然，别人访问代理获取了被缓存的响应就麻烦了。

其次，缓存失效后的重新验证也要区分开（即使用条件请求“Last-modified”和“ETag”），“must-revalidate”是只要过期就必须回源服务器验证，而新的“proxy-revalidate”只要求代理的缓存过期后必须验证，客户端不必回源，只验证到代理这个环节就行了。

再次，缓存的生存时间可以使用新的“s-maxage”（s是share的意思，注意maxage中间没有“-”），只限定在代理上能够存多久，而客户端仍然使用“max_age”。

还有一个代理专用的属性“no-transform”。代理有时候会对缓存下来的数据做一些优化，比如把图片生成png、webp等几种格式，方便今后的请求处理，而“no-transform”就会禁止这样做，不许“偷偷摸摸搞小动作”。

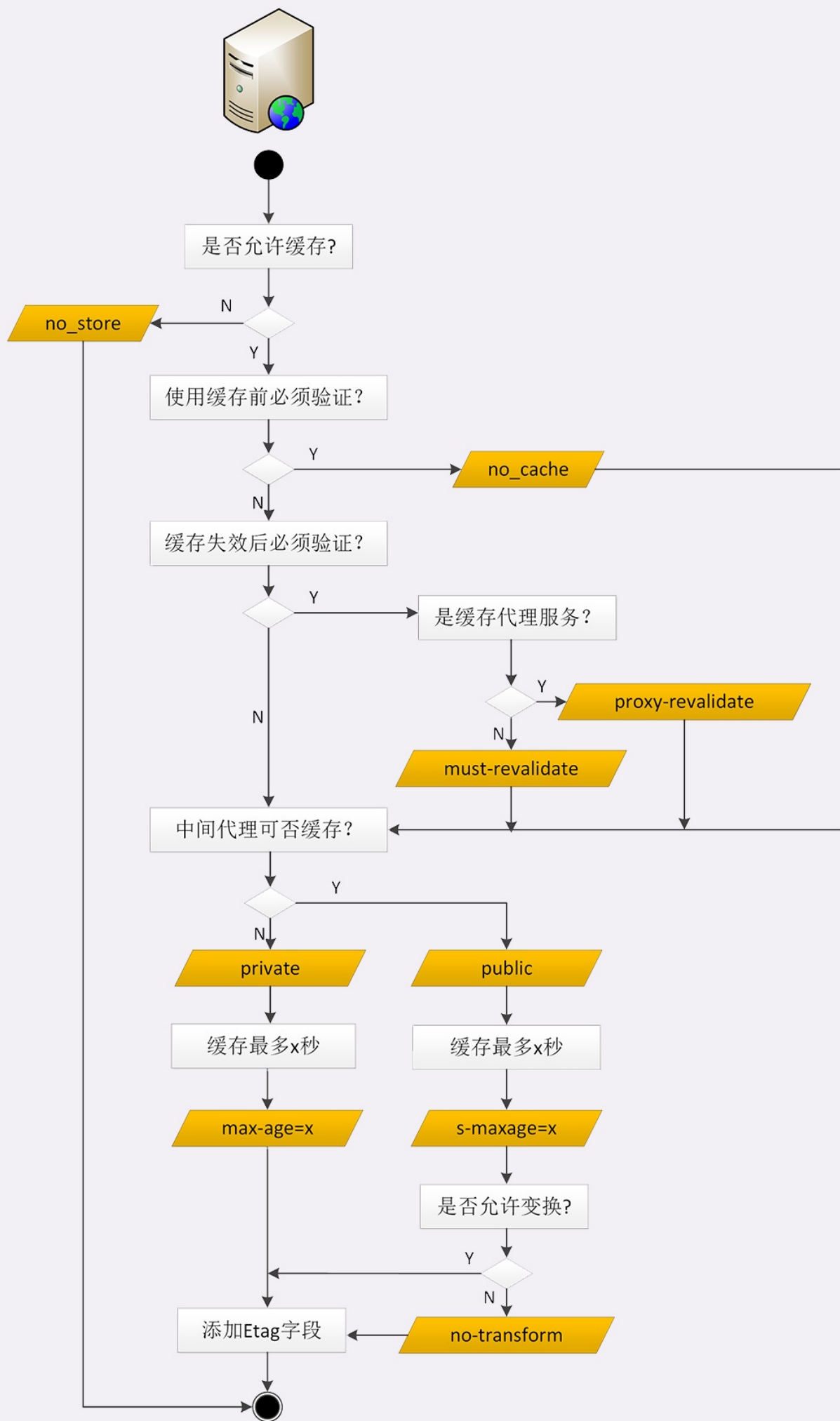
这些新的缓存控制属性比较复杂，还是用“便利店冷柜”来举例好理解一些。

水果上贴着标签“private, max-age=5”。这就是说水果不能放进冷柜，必须直接给顾客，保鲜期5天，过期了还得去超市重新进货。

冻鱼上贴着标签“public, max-age=5, s-maxage=10”。这个的意思就是可以在冰柜里存10天，但顾客那里只能存5天，过期了可以来便利店取，只要在10天之内就不必再找超市。

排骨上贴着标签“max-age=30, proxy-revalidate, no-transform”。因为缓存默认是public的，那么它在便利店和顾客的冰箱里就都可以存30天，过期后便利店必须去超市进新货，而且不能擅自把“大排”改成“小排”。

下面的流程图是完整的服务器端缓存控制策略，可以同时控制客户端和代理。

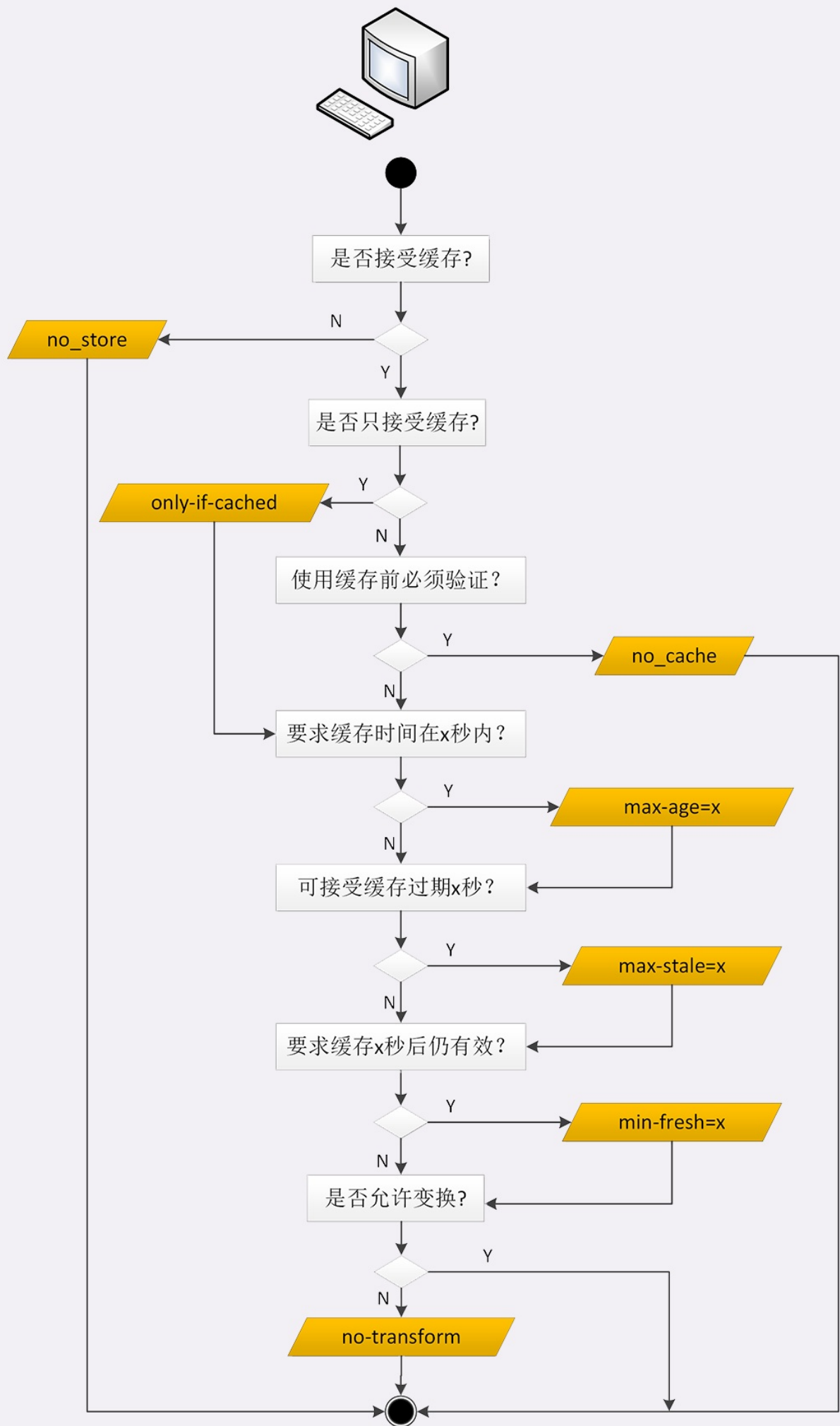


我还要提醒你一点，源服务器在设置完“Cache-Control”后必须要为报文加上“Last-modified”或“ETag”字段。否则，客户端和代理后面就无法使用条件请求来验证缓存是否有效，也就不会有304缓存重定向。

客户端的缓存控制

说完了服务器端的缓存控制策略，稍微歇一口气，我们再来看看客户端。

客户端在HTTP缓存体系里要面对的是代理和源服务器，也必须区别对待，这里我就直接上图了，来个“看图说话”。



max-age、no_store、no_cache这三个属性在[第20讲](#)已经介绍过了，它们也是同样作用于代理和源服务器。

关于缓存的生存时间，多了两个新属性“max-stale”和“min-fresh”。

“max-stale”的意思是如果代理上的缓存过期了也可以接受，但不能过期太多，超过x秒也会不要。“min-fresh”的意思是缓存必须有效，而且必须在x秒后依然有效。

比如，草莓上贴着标签“max-age=5”，现在已经在冰柜里存了7天。如果有请求“max-stale=2”，意思是过期两天也能接受，所以刚好能卖出去。

但要是“min-fresh=1”，这是绝对不允许过期的，就不会买走。这时如果有另外一个菠萝是“max-age=10”，那么“7+1<10”，在一天之后还是新鲜的，所以就能卖出去。

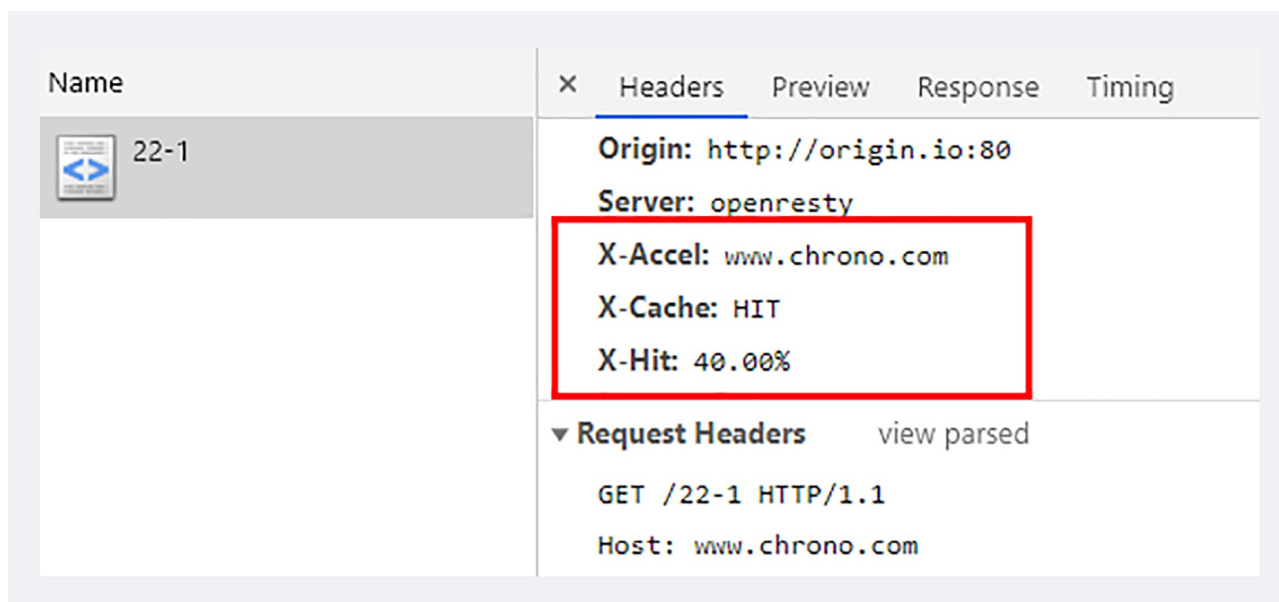
有的时候客户端还会发出一个特别的“only-if-cached”属性，表示只接受代理缓存的数据，不接受源服务器的响应。如果代理上没有缓存或者缓存过期，就应该给客户端返回一个504（Gateway Timeout）。

实验环境

信息量有些大，到这里你是不是有点头疼了，好在我们还有实验环境，用URI“/22-1”试一下吧。

它设置了“Cache-Control: public, max-age=10, s-maxage=30”，数据可以在浏览器里存10秒，在代理上存30秒，你可以反复刷新，看看代理和源服务器是怎么响应的，同样也可以配合Wireshark抓包。

代理在响应报文里还额外加了“X-Cache”“X-Hit”等自定义头字段，表示缓存是否命中和命中率，方便你观察缓存代理的工作情况。



其他问题

缓存代理的知识就快讲完了，下面再简单说两个相关的问题。

第一个是“Vary”字段，在[第15讲](#)曾经说过，它是内容协商的结果，相当于报文的一个版本标记。

同一个请求，经过内容协商后可能会有不同的字符集、编码、浏览器等版本。比如，“Vary: Accept-

Encoding” “Vary: User-Agent”，缓存代理必须要存储这些不同的版本。

当再收到相同的请求时，代理就读取缓存里的“Vary”，对比请求头里相应的“Accept-Encoding” “User-Agent”等字段，如果和上一个请求的完全匹配，比如都是“gzip” “Chrome”，就表示版本一致，可以返回缓存的数据。

另一个问题是“Purge”，也就是“缓存清理”，它对于代理也是非常重要的功能，例如：

- 过期的数据应该及时淘汰，避免占用空间；
- 源站的资源有更新，需要删除旧版本，主动换成最新版（即刷新）；
- 有时候会缓存了一些本不该存储的信息，例如网络谣言或者危险链接，必须尽快把它们删除。

清理缓存的方法有很多，比较常用的一种做法是使用自定义请求方法“PURGE”，发给代理服务器，要求删除URI对应的缓存数据。

小结

1. 计算机领域里最常用的性能优化手段是“时空转换”，也就是“时间换空间”或者“空间换时间”，HTTP缓存属于后者；
2. 缓存代理是增加了缓存功能的代理服务，缓存源服务器的数据，分发给下游的客户端；
3. “Cache-Control”字段也可以控制缓存代理，常用的有“private” “s-maxage” “no-transform”等，同样必须配合“Last-modified” “ETag”等字段才能使用；
4. 缓存代理有时候也会带来负面影响，缓存不良数据，需要及时刷新或删除。

课下作业

1. 加入了代理后HTTP的缓存复杂了很多，试着用自己的语言把这些知识再整理一下，画出有缓存代理时浏览器的工作流程图，加深理解。
2. 缓存的时间策略很重要，太大太小都不好，你觉得应该如何设置呢？

欢迎你把自己的学习体会写在留言区，与我和其他同学一起讨论。如果你觉得有所收获，也欢迎把文章分享给你的朋友。



== 课外小贴士 ==

- 01 常用的缓存代理软件有 Squid、Varnish、ATS (Apache Traffic Server) 等，而 Nginx 不仅是 Web 服务器、代理服务器，也是一个出色的缓存代理服务器，堪称“全能”。
- 02 有的缓存代理在“Cache Hit”的时候会在响应报文里加一个 Age 头字段，表示报文的生存时间，即已经在缓存里存了多久，通常它会小于“Cache-Control”里的 max-age 值，如果大于就意味着数据是“陈旧的” (stale)。
- 03 判断缓存是否命中 (Hit) 类似于查询 hash 表，使用的 key 通常就是 URI，在 Nginx 里可以用指令“proxy_cache_key”自定义。
- 04 Nginx 对 Vary 的处理实际上是做了 MD5，把 Vary 头摘要后写入缓存，请求时不仅比较 URI，也比较摘要。

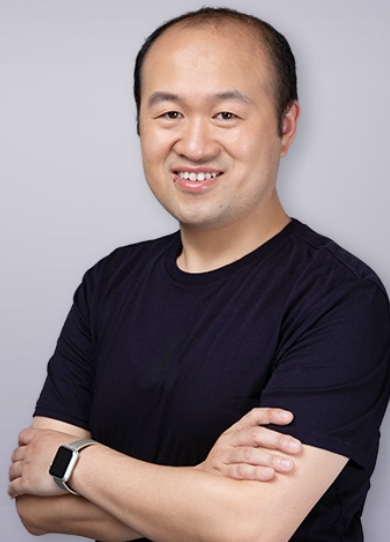
透视 HTTP 协议

深入理解 HTTP 协议本质与应用

罗剑锋

奇虎360技术专家

Nginx/OpenResty 开源项目贡献者



新版升级：点击「 请朋友读」，20位好友免费读，邀请订阅更有**现金**奖励。

精选留言：

- 院长。 2019-07-17 08:16:09

老师您好，我想请问一下，为什么有的地方说cache-control默认是private（比如cache-control的百度百科），有的地方说默认是public（比如您这篇文章），是百度百科的是错误的吗？还是根据场景不同所以默认不同吗？ [2赞]

作者回复2019-07-17 09:18:20

我看了一下rfc，对于private和public没有明确的默认值说法，可能是我弄错了，需要再测试看看。

- Geek_54edc1 2019-07-17 11:04:00

max-stale和min-fresh还是不太明白~~ [1赞]

作者回复2019-07-17 11:34:24

max-stale是可以接受的过期时间，min-fresh是可以接受的新鲜时间。

不好理解也没事，这两个属性用的不多，可以以后实际遇到了再体会。

- 居培波 2019-07-17 10:42:47

老师能结合nginx讲下缓存及代理吗？还是后面探索篇有讲。 [1赞]

作者回复2019-07-17 11:33:04

只要理解的http的缓存代理，nginx的是比较容易掌握的，可以结合nginx的文档，看proxy相关的指令，逐条对照http的功能。

单纯讲nginx就有点太大了，讲不过来。

- 何用 2019-07-17 09:24:12

老师，CDN 服务是不是就是缓存代理的一种应用？还有文中图片的 X-Accel 是什么意思呢？ [1赞]

作者回复2019-07-17 10:08:36

1，是的。

2. X-Accel是自定义字段，和x-powered-by差不多，意思是被谁加速。

- 闫飞 2019-07-18 08:13:37

min-fresh的含义是距离过期时间必须不短于约定的时间，保证取到的是短期内不会过时的内容。

但是这里有个风险是，万一服务器端因为某些原因重新刷新了资源(服务迁移等)，那么怎么反向通知缓存服务器去清理资源呢？尤其是已经返回给客户端的之前标记为fresh的资源？

作者回复2019-07-18 08:53:24

http协议里没有对此做出规定，一般的做法是由源站向代理发送pull请求，要求代理主动更新缓存。

这个pull请求不属于http协议，具体实现就看两者之间的约定了。

- 一步 2019-07-17 20:40:56

老师请教一个有关302的问题，就是再前端代码中有个post请求，去请求后端服务器，当后端服务器处理完业务逻辑后，需要重定向到另一个网站，返回的是一个302的状态码和响应头Location是另一个网站地址。这时候问题出现了，当返回到前端的时候，浏览器没有自动跳转重定向后的地址，而是当作接口去请求了需要跳转后的地址，然后就出现了跨域的问题

这个是什么原因呢？

作者回复2019-07-18 08:54:55

可以参考一下第18讲，改用303 see other。

- WL 2019-07-17 19:26:58

请问一下老师头摘要Vary也是一个header字段吗？这个具体怎么用，以前写项目好像没有用到过

作者回复2019-07-17 20:07:53

vary主要用在缓存，用来告诉缓存代理此报文依据的是哪些请求头字段。

缓存代理收到带vary的响应报文后，就会根据vary里的字段提取请求头，计算hash，然后和缓存存在一起。

当有新请求的时候就比较新请求的头字段hash与缓存的hash，一致就说明可以重用缓存。

- 徐海浪 2019-07-17 13:14:55

我们的做法是源服务器跟代理服务器联动，源服务器有文件变化(版本发布)，触发更新代理服务器缓存。如果没有联动的机制，简单粗暴根据应用升级周期设置过期时间也可以，但这样会多做无用功。

作者回复2019-07-17 14:43:12

欢迎经验分享！