

你还在为 TCP 重传、滑动窗口、流量控制、拥塞控制发愁吗？看完图解就不愁了

[小林coding](#)

每日一句英语学习，每天进步一点点：

前言

前一篇「[硬不硬你说了算！近 40 张图解被问千百遍的 TCP 三次握手和四次挥手面试题](#)」得到了很多读者的认可，在此特别感谢你们的认可，大家都暖暖的。



突然害羞

来了，今天又来图解 TCP 了，小林可能会迟到，但不会缺席。

迟到的原因，主要是 TCP 巨复杂，它为了保证可靠性，用了巨多的机制来保证，真是个「伟大」的协议，写着写着发现这水太深了。。。

本文的全部图片都是小林绘画的，非常的辛苦且累，不废话了，直接进入正文吧，Go!

正文

相信大家都知道 TCP 是一个可靠传输的协议，那它是如何保证可靠的呢？

为了实现可靠性传输，需要考虑很多事情，例如数据的破坏、丢包、重复以及分片顺序混乱等问题。如不能解决这些问题，也就无从谈起可靠传输。

那么，TCP 是通过序列号、确认应答、重发控制、连接管理以及窗口控制等机制实现可靠性传输的。

今天，将重点介绍 TCP 的重传机制、滑动窗口、流量控制、拥塞控制。

提纲

重传机制

TCP 实现可靠传输的方式之一，是通过序列号与确认应答。

在 TCP 中，当发送端的数据到达接收主机时，接收端主机会返回一个确认应答消息，表示已收到消息。

但在错综复杂的网络，并不一定能如上图那么顺利能正常的数据传输，万一数据在传输过程中丢失了呢？

所以 TCP 针对数据包丢失的情况，会用重传机制解决。

接下来说说常见的重传机制：

- 超时重传
- 快速重传
- SACK
- D-SACK

超时重传

重传机制的其中一个方式，就是在发送数据时，设定一个定时器，当超过指定的时间后，没有收到对方的 ACK 确认应答报文，就会重发该数据，也就是我们常说的**超时重传**。

TCP 会在以下两种情况发生超时重传：

- 数据包丢失
- 确认应答丢失

超时重传的两种情况

超时时间应该设置为多少呢？

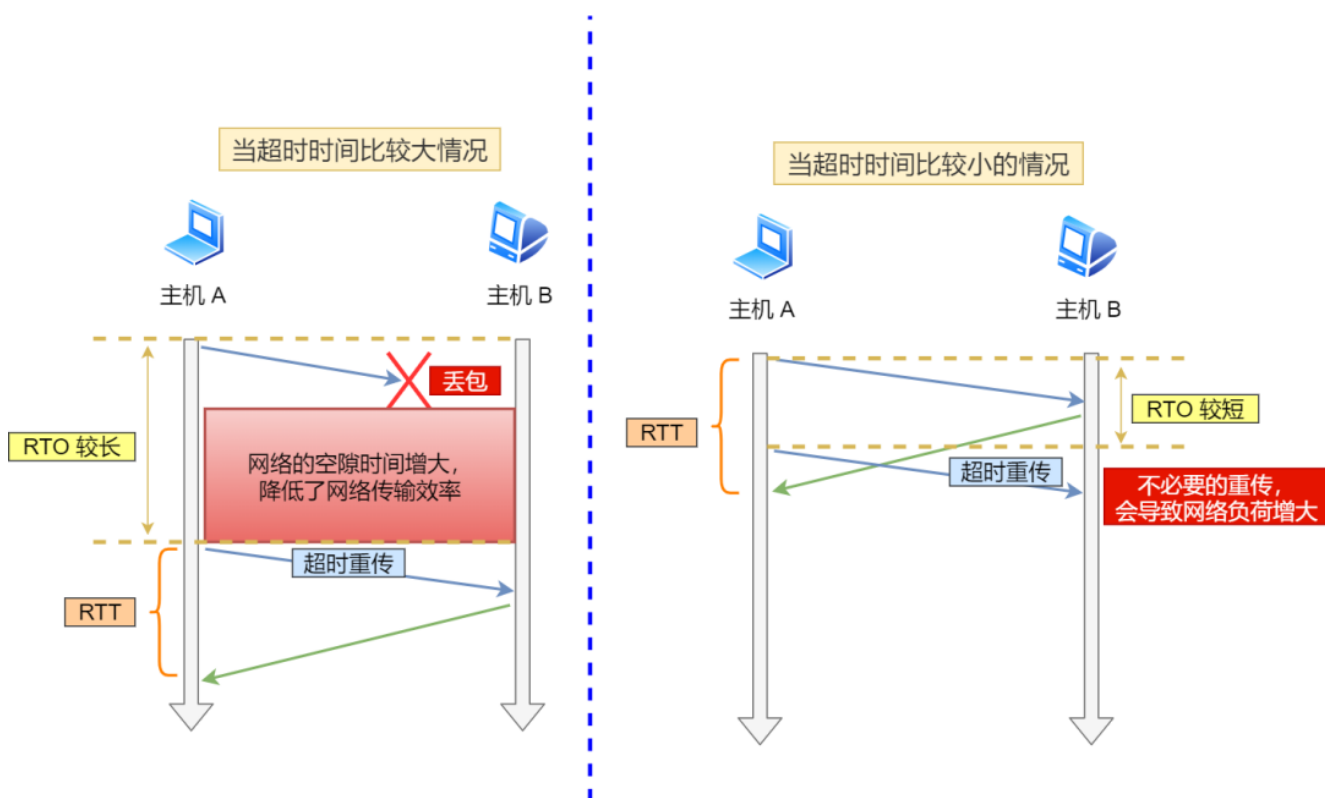
我们先来了解一下什么是 RTT (Round-Trip Time 往返时延)，从下图我们就可以知道：

RTT

RTT 就是数据从网络一端传送到另一端所需的时间，也就是包的往返时间。

超时重传时间是以 RTO (Retransmission Timeout 超时重传时间) 表示。

假设在重传的情况下，超时时间 RTO 「较长或较短」时，会发生什么事情呢？



超时时间较长与较短

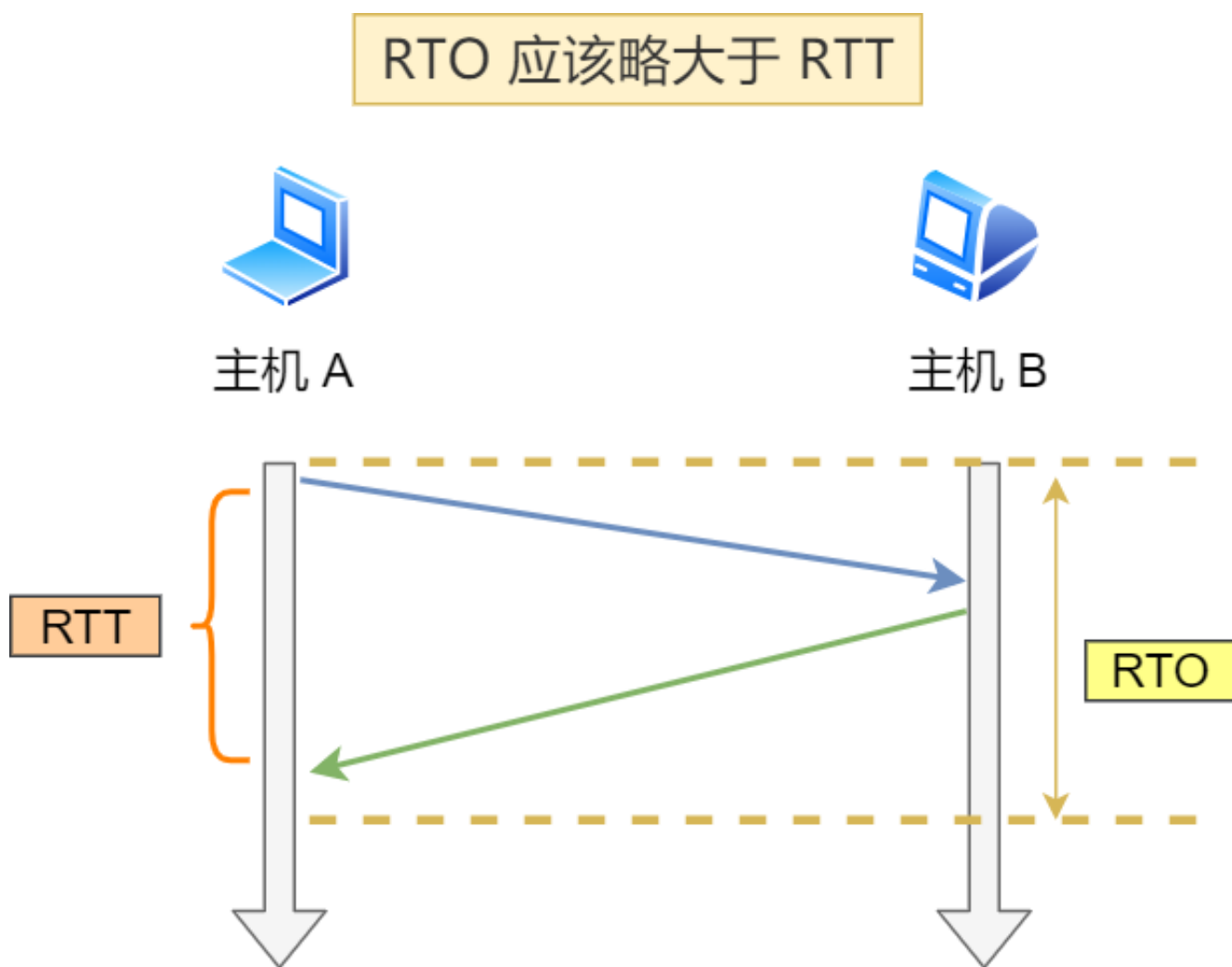
上图中有两种超时时间不同的情况：

- 当超时时间 **RTO** 较大时，重发就慢，丢了老半天才重发，没有效率，性能差；

- 当超时时间 **RTO** 较小时，会导致可能并没有丢就重发，于是重发的就快，会增加网络拥塞，导致更多的超时，更多的超时导致更多的重发。

精确的测量超时时间 RTO 的值是非常重要的，这可让我们的重传机制更高效。

根据上述的两种情况，我们可以得知，超时重传时间 **RTO** 的值应该略大于报文往返 **RTT** 的值。



RTO 应略大于 RTT

至此，可能大家觉得超时重传时间 RTO 的值计算，也不是很复杂嘛。

好像就是在发送端发包时记下 t_0 ，然后接收端再把这个

ack 回来时再记一个 t_1 ，于是 $RTT = t_1 - t_0$ 。没那么简单，这只是一个采样，不能代表普遍情况。

实际上「报文往返 RTT 的值」是经常变化的，因为我们的网络也是时常变化的。也就因为「报文往返 RTT 的值」是经常波动变化的，所以「超时重传时间 RTO 的值」应该是一个动态变化的值。

我们来看看 Linux 是如何计算 RTO 的呢？

估计往返时间，通常需要采样以下两个：

- 需要 TCP 通过采样 RTT 的时间，然后进行加权平均，算出一个平滑 RTT 的值，而且这个值还是要不断变化的，因为网络状况不断地变化。
- 除了采样 RTT，还要采样 RTT 的波动范围，这样就避免如果 RTT 有一个大的波动的话，很难被发现的情况。

RFC6289 建议使用以下的公式计算 RTO：

① 首次计算 RTO，其中 R1 为第一次测量的 RTT

$$SRTT = R1$$

$$DevRTT = R1/2$$

$$RTO = \mu * SRTT + \partial * DevRTT = \mu * R1 + \partial * (R1/2)$$

② 后续计算 RTO，其中 R2 为最新测量的 RTT

$$SRTT = SRTT + \alpha (RTT - SRTT) = R1 + \alpha * (R2 - R1)$$

$$DevRTT = (1-\beta) * DevRTT + \beta * (|RTT - SRTT|) = (1-\beta) * (R1/2) + \beta * (|R2 - R1|)$$

$$RTO = \mu * SRTT + \partial * DevRTT$$

RFC6289 建议的 RTO 计算

其中 SRTT 是计算平滑的 RTT，DevRTR 是计算平滑的 RTT 与最新 RTT 的差距。

在 Linux 下， $\alpha = 0.125$ ， $\beta = 0.25$ ， $\mu = 1$ ， $\partial = 4$ 。别问怎么来的，问就是大量实验中调出来的。

如果超时重发的数据，再次超时的时候，又需要重传的时候，TCP 的策略是超时间隔加倍。

也就是每当遇到一次超时重传的时候，都会将下一次超时时间间隔设为先前值的两倍。两次超时，就说明网络环境差，不宜频繁反复发送。

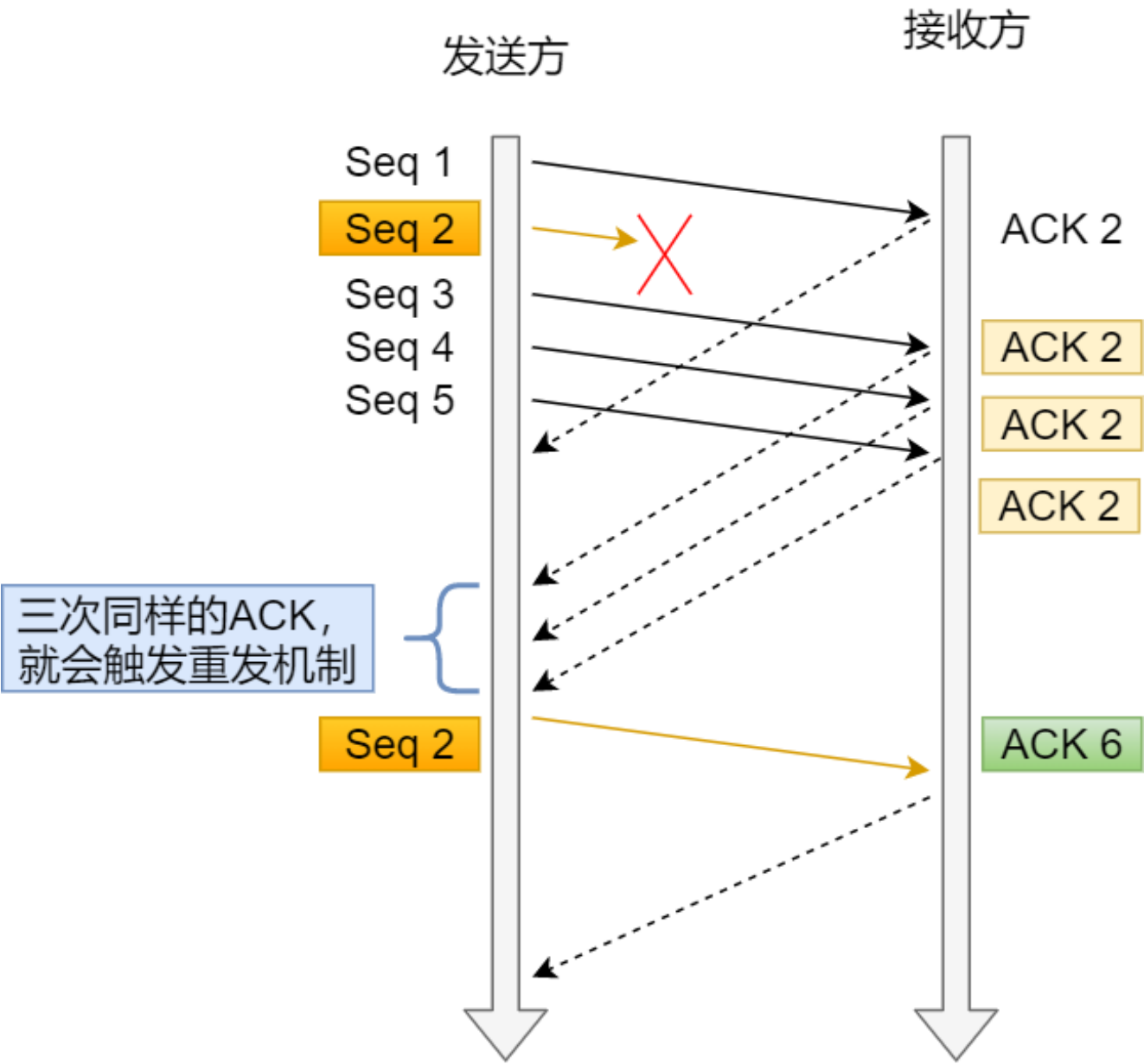
超时触发重传存在的问题是，超时周期可能相对较长。那是不是可以有更快的方式呢？

于是就可以用「快速重传」机制来解决超时重发的时间等待。

快速重传

TCP 还有另外一种快速重传（**Fast Retransmit**）机制，它不以时间为驱动，而是以数据驱动重传。

快速重传机制，是如何工作的呢？其实很简单，一图胜千言。



快速重传机制

在上图，发送方发出了 1，2，3，4，5 份数据：

- 第一份 Seq1 先送到了，于是就 Ack 回 2；
- 结果 Seq2 因为某些原因没收到，Seq3 到达了，于是还是 Ack 回 2；
- 后面的 Seq4 和 Seq5 都到了，但还是 Ack 回 2，因为 Seq2 还是没有收到；
- 发送端收到了三个 **Ack = 2** 的确认，知道了 **Seq2** 还没有收到，就会在定时器过期之前，重传丢失的 **Seq2**。
- 最后，接收到收到了 Seq2，此时因为 Seq3，Seq4，Seq5 都收到了，于是 Ack 回 6。

所以，快速重传的工作方式是当收到三个相同的 ACK 报文时，会在定时器过期之前，重传丢失的报文段。

快速重传机制只解决了一个问题，就是超时时间的问题，但是它依然面临着另外一个问题。就是重传的时候，是重传之前的一个，还是重传所有的问题。

比如对于上面的例子，是重传 Seq2 呢？还是重传 Seq2、Seq3、Seq4、Seq5 呢？因为发送端并不清楚这连续的三个 Ack 2 是谁传回来的。

根据 TCP 不同的实现，以上两种情况都是有可能的。可见，这是一把双刃剑。

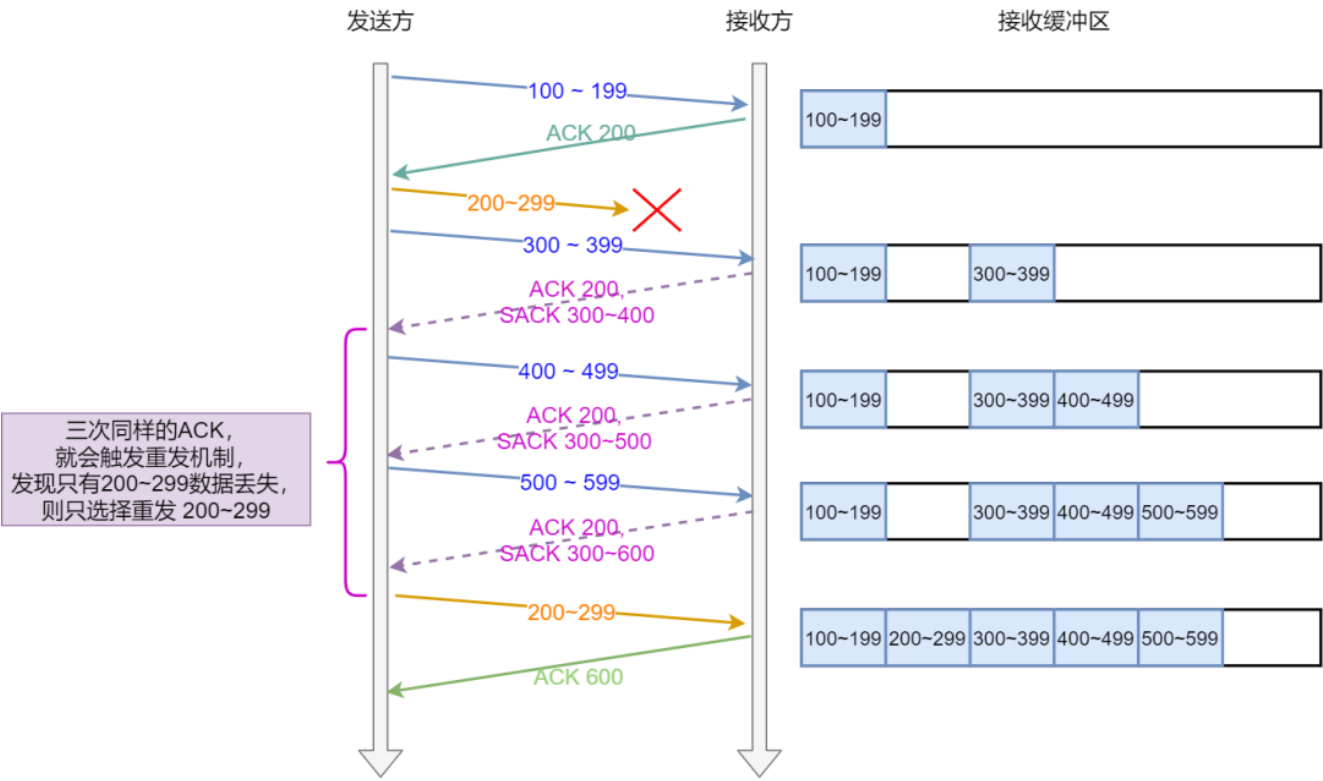
为了解决不知道该重传哪些 TCP 报文，于是就有 SACK 方法。

SACK 方法

还有一种实现重传机制的方式叫：sack（Selective Acknowledgment 选择性确认）。

这种方式需要在 TCP 头部「选项」字段里加一个 sack 的东西，它可以将缓存的地图发送给发送方，这样发送方就可以知道哪些数据收到了，哪些数据没收到，知道了这些信息，就可以只重传丢失的数据。

如下图，发送方收到了三次同样的 ACK 确认报文，于是就会触发快速重发机制，通过 sack 信息发现只有 200~299 这段数据丢失，则重发时，就只选择了这个 TCP 段进行重复。



选择性确认

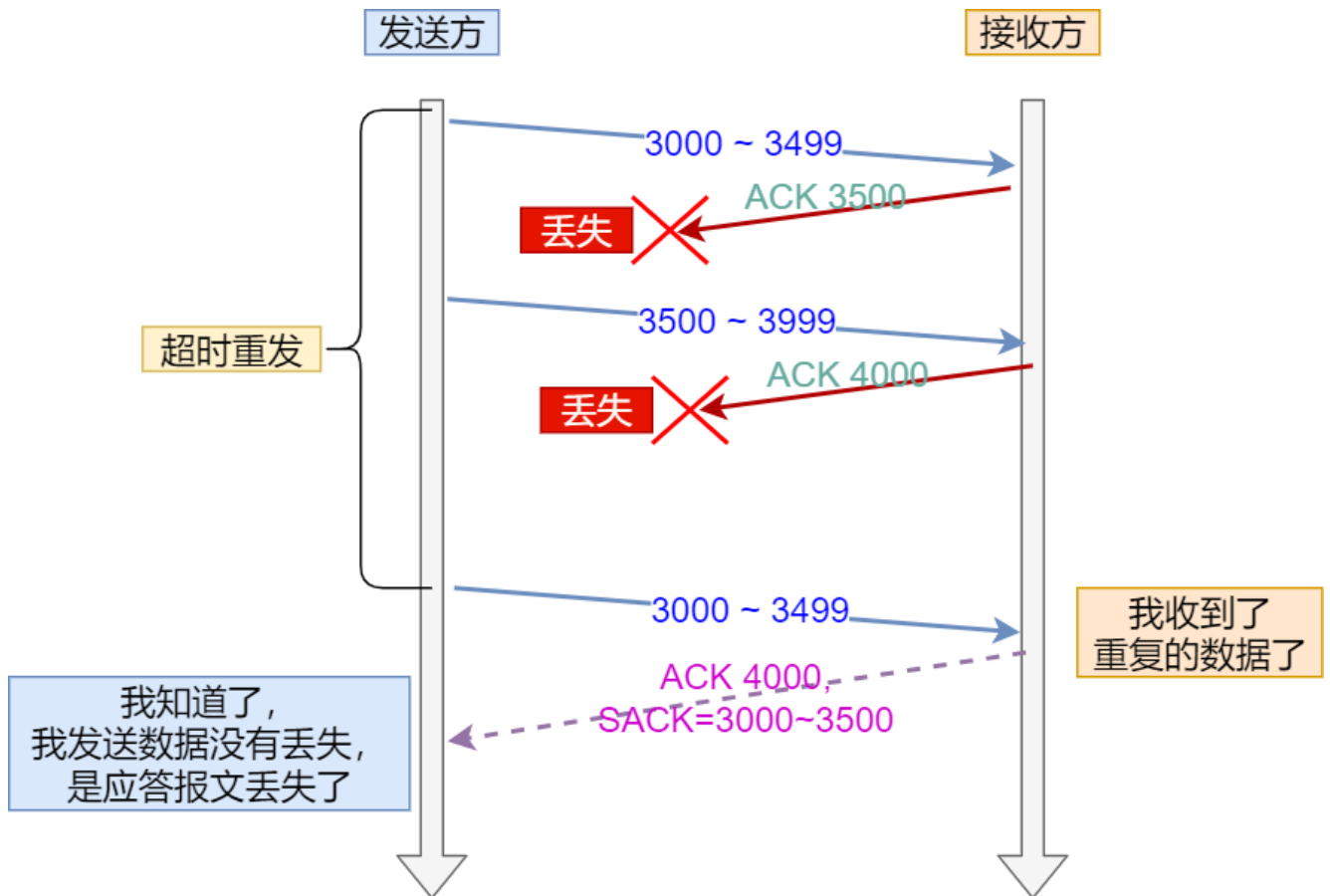
如果要支持 sack，必须双方都要支持。在 Linux 下，可以通过 `net.ipv4.tcp_sack` 参数打开这个功能（Linux 2.4 后默认打开）。

Duplicate SACK

Duplicate SACK 又称 D-SACK，其主要使用了 **SACK** 来告诉「发送方」有哪些数据被重复接收了。

下面举例两个栗子，来说明 D-SACK 的作用。

栗子一号：ACK 丢包



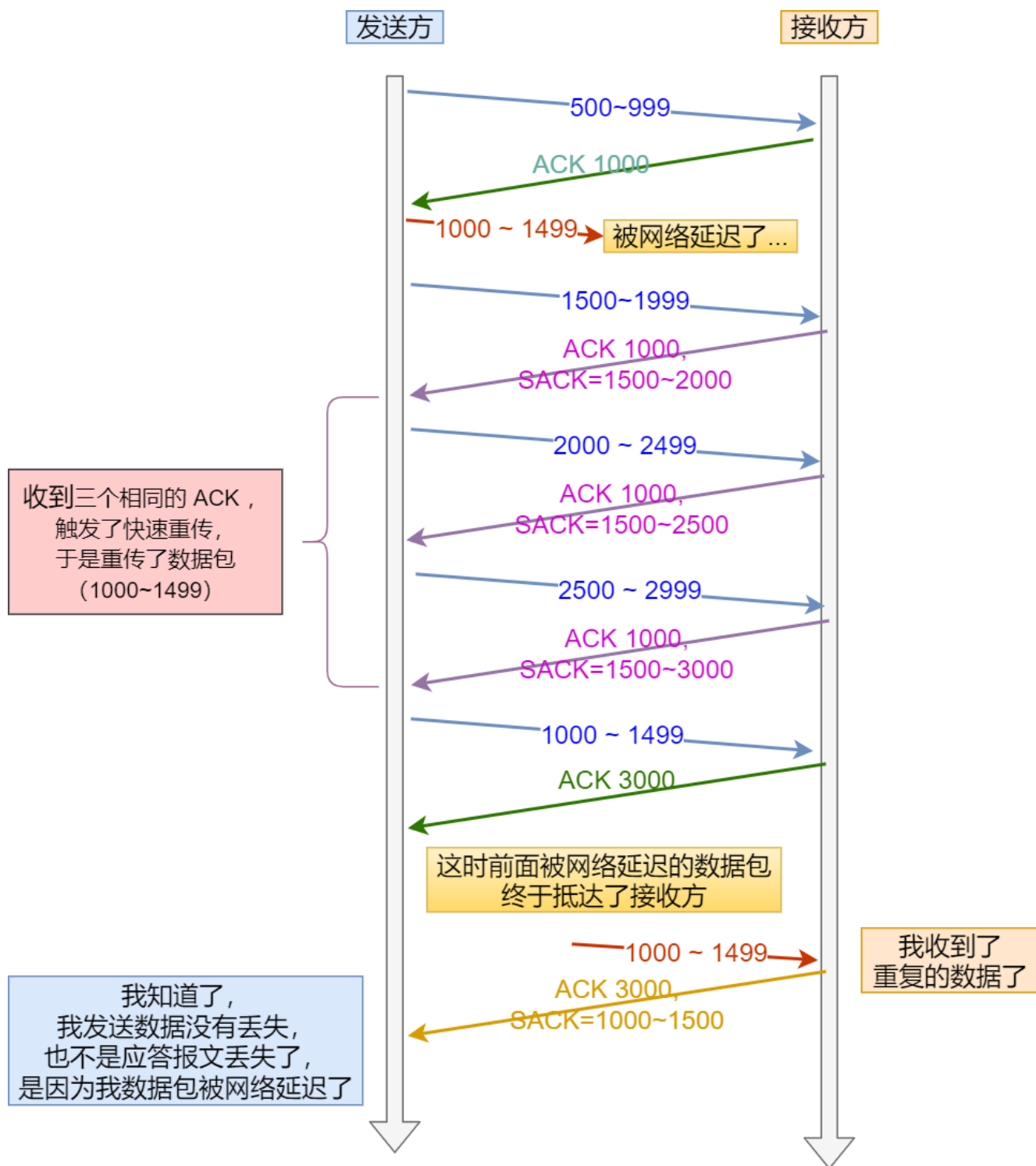
ACK 丢包

- 「接收方」发给「发送方」的两个 ACK 确认应答都丢失了，所以发送方超时后，重传第一个数据包 (3000 ~ 3499)
- 于是「接收方」发现数据是重复收到的，于是回了一个 **SACK = 3000~3500**，告诉「发送方」3000~3500 的数据早已被接收了，因为 ACK 都到了 4000 了，已经意味着 4000 之前的所有数据都已收

到，所以这个 SACK 就代表着 D-SACK。

- 这样「发送方」就知道了，数据没有丢，是「接收方」的 ACK 确认报文丢了。

栗子二号：网络延时



网络延时

- 数据包（1000~1499）被网络延迟了，导致「发送

方」没有收到 Ack 1500 的确认报文。

- 而后面报文到达的三个相同的 ACK 确认报文，就触发了快速重传机制，但是在重传后，被延迟的数据包（1000~1499）又到了「接收方」；
- 所以「接收方」回了一个 **SACK=1000~1500**，因为 **ACK** 已经到了 **3000**，所以这个 **SACK** 是 **D-SACK**，表示收到了重复的包。
- 这样发送方就知道快速重传触发的原因不是发出去的包丢了，也不是因为回应的 ACK 包丢了，而是因为网络延迟了。

可见，D-SACK 有这么几个好处：

1. 可以让「发送方」知道，是发出去的包丢了，还是接收方回应的 ACK 包丢了；
2. 可以知道是不是「发送方」的数据包被网络延迟了；
3. 可以知道网络中是不是把「发送方」的数据包给复制了；

在 Linux 下可以通过 `net.ipv4.tcp_dsack` 参数开启/关闭这个功能（Linux 2.4 后默认打开）。

滑动窗口

引入窗口概念的原因

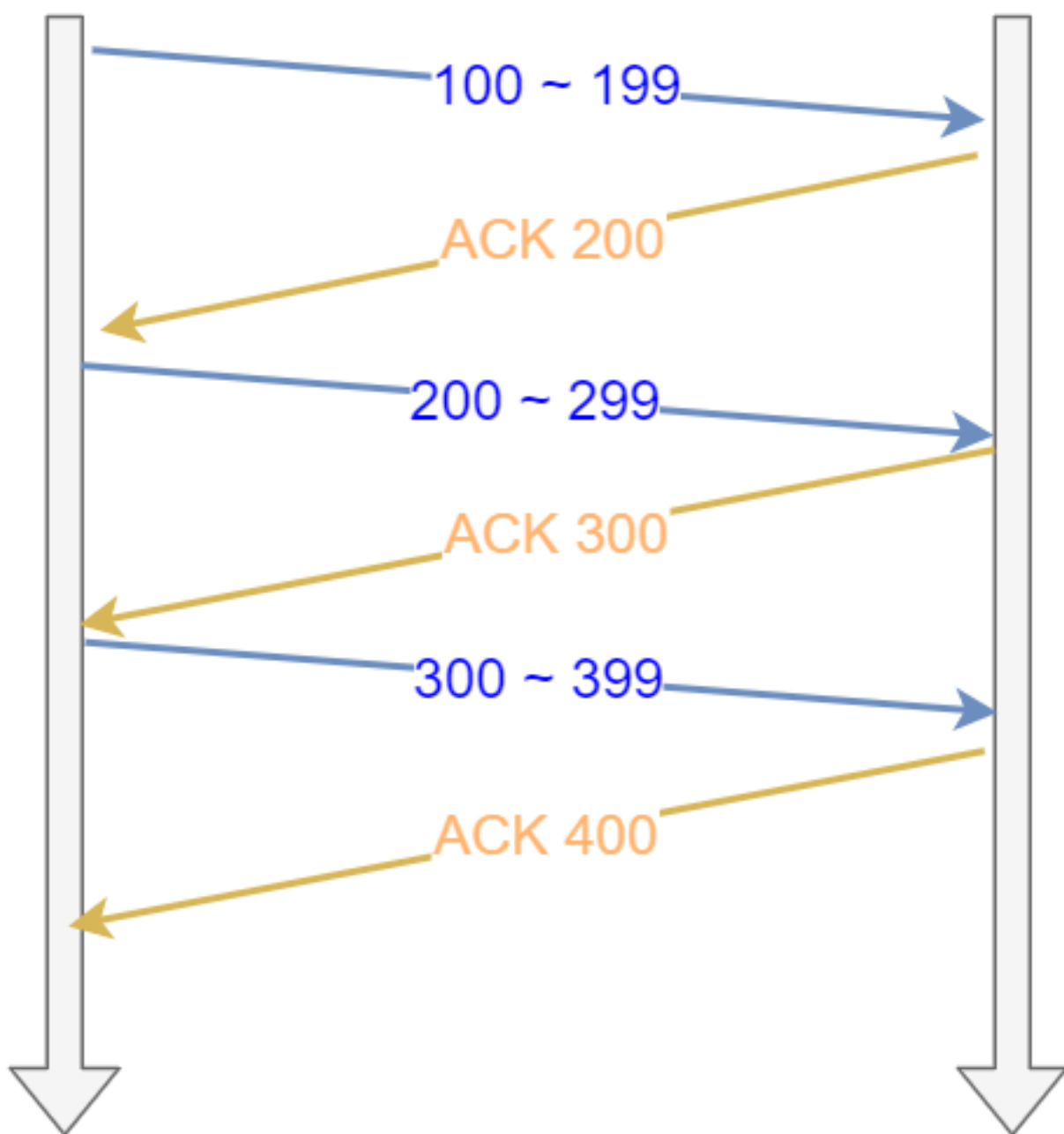
我们都知道 TCP 是每发送一个数据，都要进行一次确认应答。当上一个数据包收到了应答了，再发送下一个。

这个模式就有点像我和你面对面聊天，你一句我一句。但这种方式的缺点是效率比较低的。

如果你说完一句话，我在处理其他事情，没有及时回复你，那你不是要干等着我做完了其他事情后，我回复你，你才能说下一句话，很显然这不现实。

发送方

接收方



为每个数据包确认应答的缺点：
包的往返时间越长，网络的吞吐量会越低

按数据包进行确认应答

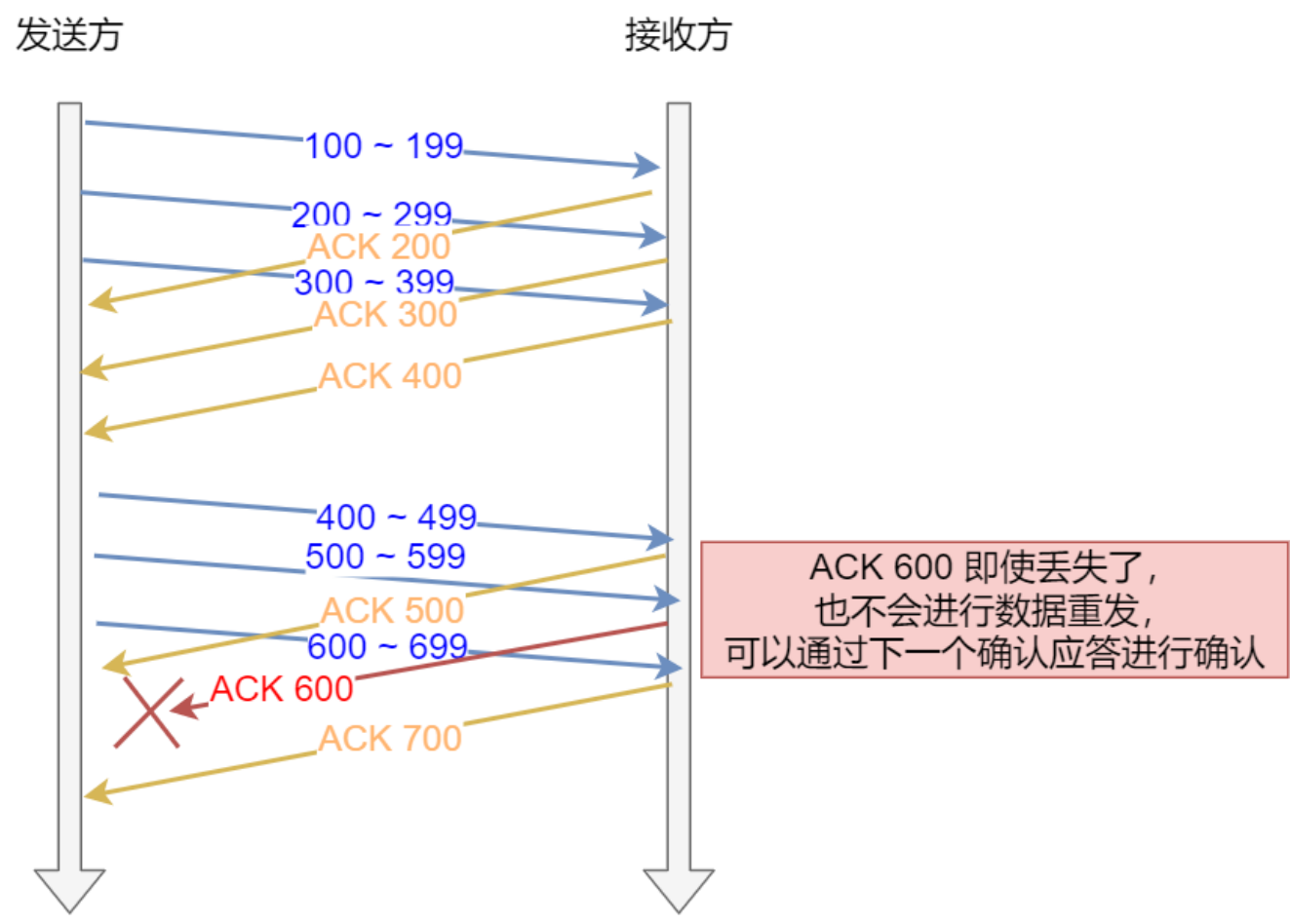
所以，这样的传输方式有一个缺点：数据包的往返时间越长，通信的效率就越低。

为了解决这个问题，TCP 引入了窗口这个概念。即使在往返时间较长的情况下，它也不会降低网络通信的效率。

那么有了窗口，就可以指定窗口大小，窗口大小就是指无需等待确认应答，而可以继续发送数据的最大值。

窗口的实现实际上是操作系统开辟的一个缓存空间，发送方主机在等到确认应答返回之前，必须在缓冲区中保留已发送的数据。如果按期收到确认应答，此时数据就可以从缓存区清除。

假设窗口大小为 3 个 TCP 段，那么发送方就可以「连续发送」 3 个 TCP 段，并且中途若有 ACK 丢失，可以通过「下一个确认应答进行确认」。如下图：



用滑动窗口方式并行处理

图中的 ACK 600 确认应答报文丢失，也没关系，因为可以通下一个确认应答进行确认，只要发送方收到了 ACK 700 确认应答，就意味着 700 之前的所有数据「接收方」都收到了。这个模式就叫累计确认或者累计应答。

窗口大小由哪一方决定？

TCP 头里有一个字段叫 window，也就是窗口大小。

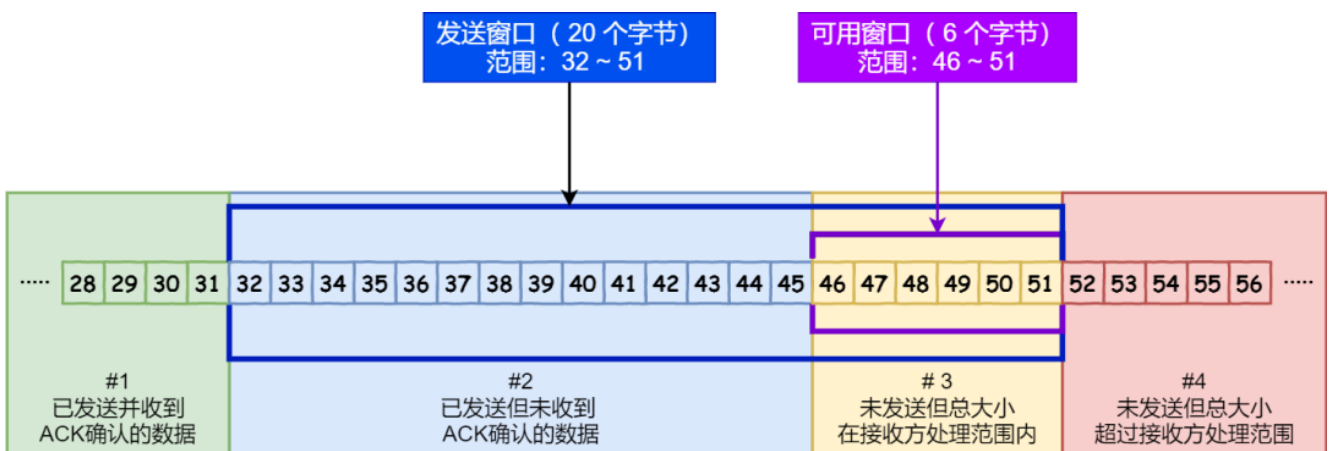
这个字段是接收端告诉发送端自己还有多少缓冲区可以接收数据。于是发送端就可以根据这个接收端的处理能力来发送数据，而不会导致接收端处理不过来。

所以，通常窗口的大小是由接收方的决定的。

发送方发送的数据大小不能超过接收方的窗口大小，否则接收方就无法正常接收到数据。

发送方的滑动窗口

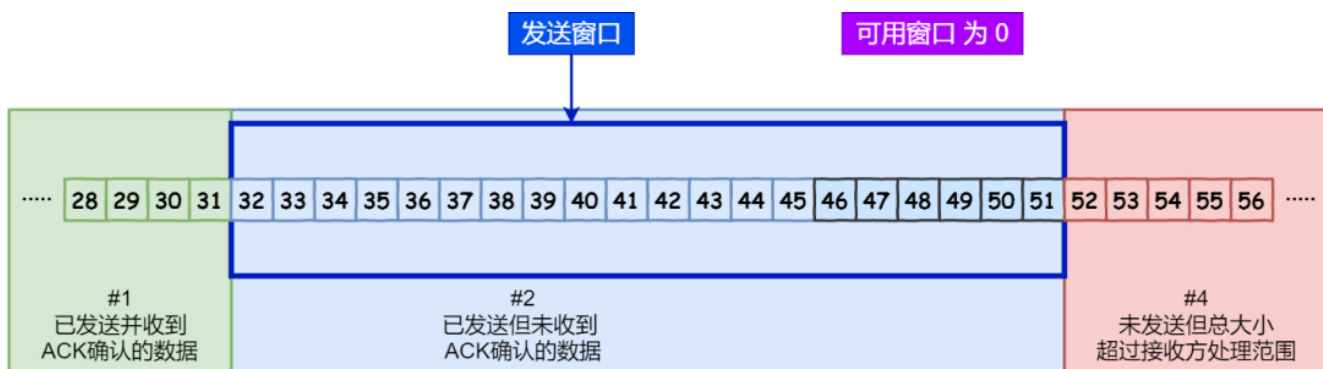
我们先来看看发送方的窗口，下图就是发送方缓存的数据，根据处理的情况分成四个部分，其中深蓝色方框是发送窗口，紫色方框是可用窗口：



- #1 是已发送并收到 ACK 确认的数据：1~31 字节
- #2 是已发送但未收到 ACK 确认的数据：32~45 字节
- #3 是未发送但总大小在接收方处理范围内（接收方还有空间）：46~51 字节

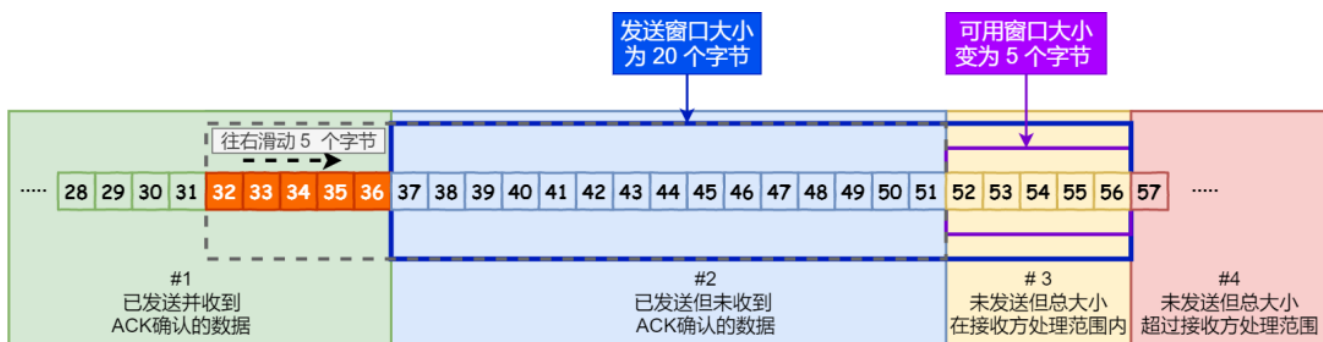
- #4 是未发送但总大小超过接收方处理范围（接收方没有空间）：52字节以后

在下图，当发送方把数据「全部」都一下发送出去后，可用窗口的大小就为 0 了，表明可用窗口耗尽，在没收到 ACK 确认之前是无法继续发送数据了。



可用窗口耗尽

在下图，当收到之前发送的数据 32~36 字节的 ACK 确认应答后，如果发送窗口的大小没有变化，则滑动窗口往右边移动 5 个字节，因为有 5 个字节的数据被应答确认，接下来 52~56 字节又变成了可用窗口，那么后续也就可以发送 52~56 这 5 个字节的数据了。

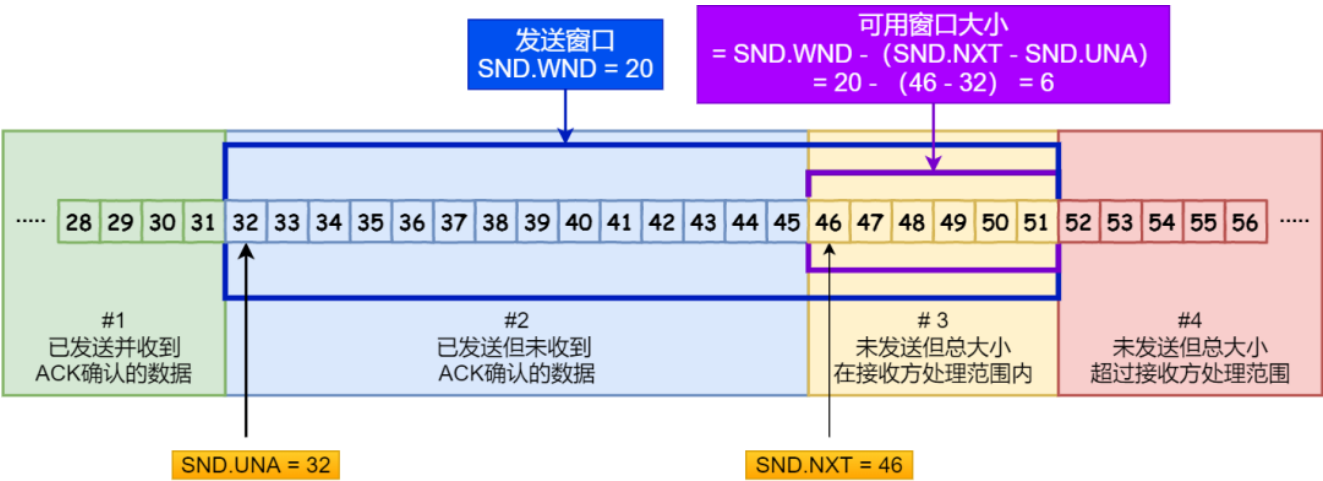


32 ~ 36 字节已确认

程序是如何表示发送方的四个部分的呢？

TCP 滑动窗口方案使用三个指针来跟踪在四个传输类别中的每一个类别中的字节。其中两个指针是绝对指针（指特

定的序列号)，一个是相对指针（需要做偏移）。



SND.WND、SND.UN、SND.NXT

- `SND.WND`：表示发送窗口的大小（大小是由接收方指定的）；
- `SND.UNA`：是一个绝对指针，它指向的是已发送但未收到确认的第一个字节的序列号，也就是 #2 的第一个字节。
- `SND.NXT`：也是一个绝对指针，它指向未发送但可发送范围的第一个字节的序列号，也就是 #3 的第一个字节。
- 指向 #4 的第一个字节是个相对指针，它需要 `SND.UNA` 指针加上 `SND.WND` 大小的偏移量，就可以指向 #4 的第一个字节了。

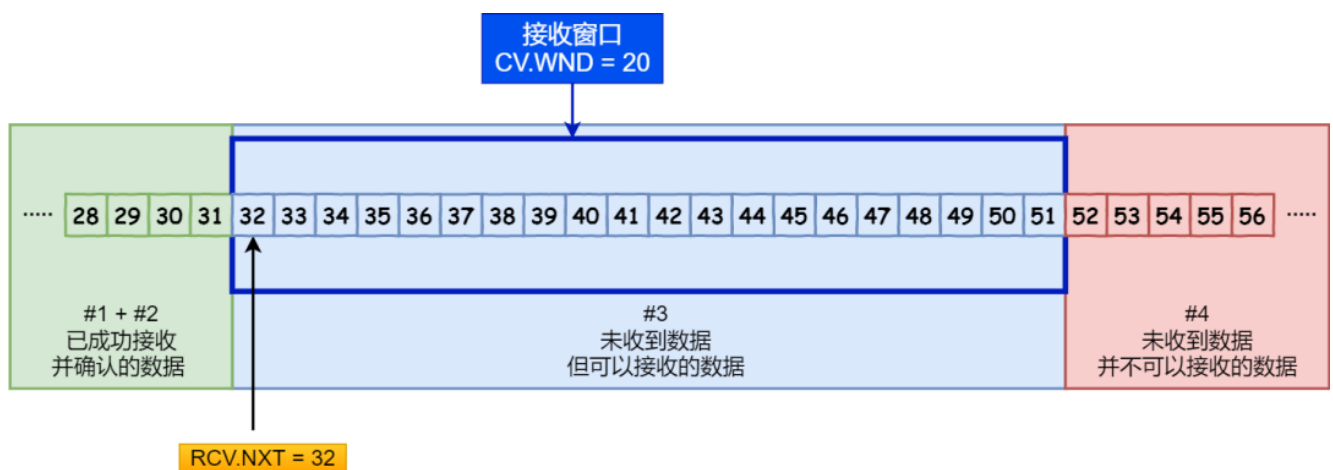
那么可用窗口大小的计算就可以是：

可用窗口大 = **SND.WND - (SND.NXT - SND.UNA)**

接收方的滑动窗口

接下来我们看看接收方的窗口，接收窗口相对简单一些，根据处理的情况划分成三个部分：

- #1 + #2 是已成功接收并确认的数据（等待应用进程读取）；
- #3 是未收到数据但可以接收的数据；
- #4 未收到数据并不可以接收的数据；



接收窗口

其中三个接收部分，使用两个指针进行划分：

- `RCV.WND`：表示接收窗口的大小，它会通告给发送方。
- `RCV.NXT`：是一个指针，它指向期望从发送方发送来的下一个数据字节的序列号，也就是 #3 的第一个字节。
- 指向 #4 的第一个字节是个相对指针，它需要 `RCV.NXT` 指针加上 `RCV.WND` 大小的偏移量，就可以指向 #4 的第一个字节了。

接收窗口和发送窗口的大小是相等的吗？

并不是完全相等，接收窗口的大小是约等于发送窗口的大小的。

因为滑动窗口并不是一成不变的。比如，当接收方的应用进程读取数据的速度非常快的话，这样的话接收窗口可以很快的就空缺出来。那么新的接收窗口大小，是通过 TCP 报文中的 Windows 字段来告诉发送方。那么这个传输过程是存在时延的，所以接收窗口和发送窗口是约等于的关系。

流量控制

发送方不能无脑的发数据给接收方，要考虑接收方处理能力。

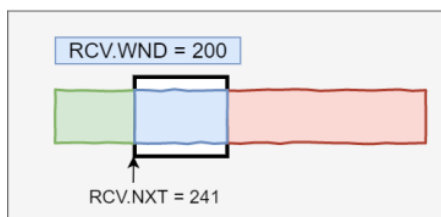
如果一直无脑的发数据给对方，但对方处理不过来，那么就会导致触发重发机制，从而导致网络流量的无端的浪费。

为了解决这种现象发生，**TCP** 提供一种机制可以让「发送方」根据「接收方」的实际接收能力控制发送的数据量，这就是所谓的流量控制。

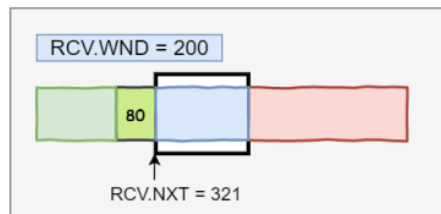
下面举个栗子，为了简单起见，假设以下场景：

- 客户端是接收方，服务端是发送方
- 假设接收窗口和发送窗口相同，都为 200

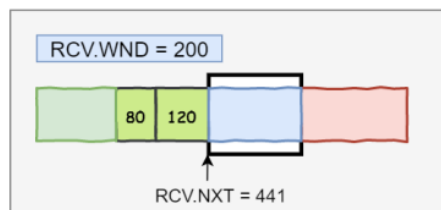
- 假设两个设备在整个传输过程中都保持相同的窗口大小，不受外界影响



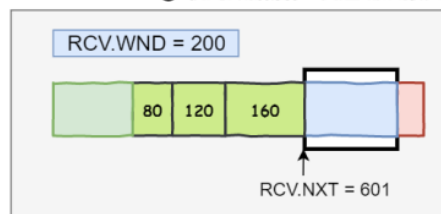
① 发送请求数据



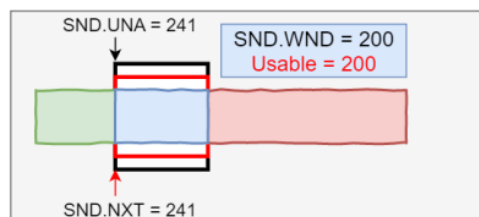
③ 收到数据后，接着发送确认报文



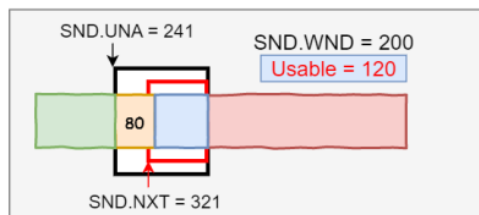
⑤ 收到数据后，发送确认报文



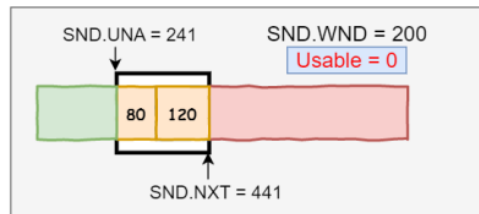
⑨ 收到数据后，发送确认报文



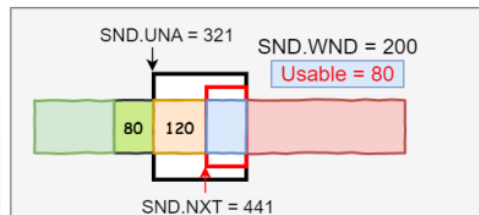
② 收到请求，发送确认和 80 字节的数据



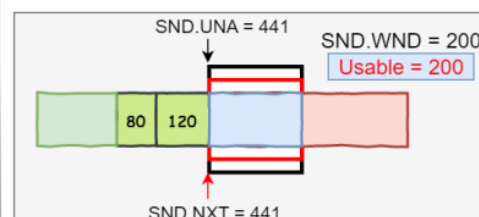
④ 发送 120 个字节数据



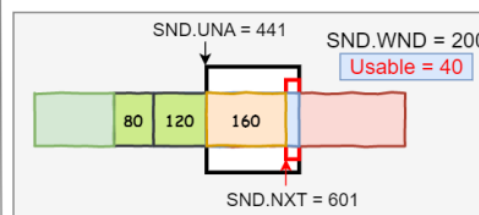
⑥ 收到确认报文后，窗口往右移动 80 字节



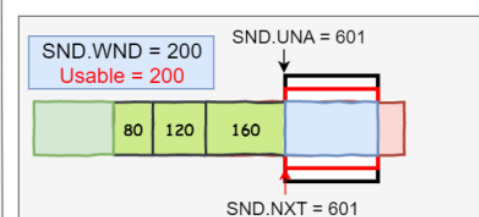
⑦ 收到确认报文后，窗口往右移动 200 字节



⑧ 发送 160 个字节数据



⑩ 收到确认报文后，窗口往右移动 160 字节

Length 140
Seq num = 1Length = 80
Seq Num = 241
Ack Num = 141

Ack Num = 321

Length = 120
Seq Num = 321

Ack Num = 441

Length = 160
Seq num = 441

Ack Num = 601

根据上图的流量控制，说明下每个过程：

1. 客户端向服务端发送请求数据报文。这里要说明下，本次例子是把服务端作为发送方，所以没有画出服务端的接收窗口。
2. 服务端收到请求报文后，发送确认报文和 80 字节的数据，于是可用窗口 `Usable` 减少为 120 字节，同时 `SND.NXT` 指针也向右偏移 80 字节后，指向 321，这意味着下次发送数据的时候，序列号是 **321**。
3. 客户端收到 80 字节数据后，于是接收窗口往右移动 80 字节，`RCV.NXT` 也就指向 321，这意味着客户端期望的下一个报文的序列号是 **321**，接着发送确认报文给服务端。
4. 服务端再次发送了 120 字节数据，于是可用窗口耗尽为 0，服务端无法再继续发送数据。
5. 客户端收到 120 字节的数据后，于是接收窗口往右移动 120 字节，`RCV.NXT` 也就指向 441，接着发送确认报文给服务端。
6. 服务端收到对 80 字节数据的确认报文后，`SND.UNA` 指针往右偏移后指向 321，于是可用窗口 `Usable` 增大到 80。
7. 服务端收到对 120 字节数据的确认报文后，`SND.UNA` 指针往右偏移后指向 441，于是可用窗口 `Usable` 增

大到 200。

8. 服务端可以继续发送了，于是发送了 160 字节的数据后，`SND.NXT` 指向 601，于是可用窗口 `Usable` 减少到 40。
9. 客户端收到 160 字节后，接收窗口往右移动了 160 字节，`RCV.NXT` 也就是指向了 601，接着发送确认报文给服务端。
10. 服务端收到对 160 字节数据的确认报文后，发送窗口往右移动了 160 字节，于是 `SND.UNA` 指针偏移了 160 后指向 601，可用窗口 `Usable` 也就增大至了 200。

操作系统缓冲区与滑动窗口的关系

前面的流量控制例子，我们假定了发送窗口和接收窗口是不变的，但是实际上，发送窗口和接收窗口中所存放的字节数，都是放在操作系统内存缓冲区中的，而操作系统的缓冲区，会被操作系统调整。

当应用进程没办法及时读取缓冲区的内容时，也会对我们的缓冲区造成影响。

那操心系统的缓冲区，是如何影响发送窗口和接收窗口的呢？

我们先来看看第一个例子。

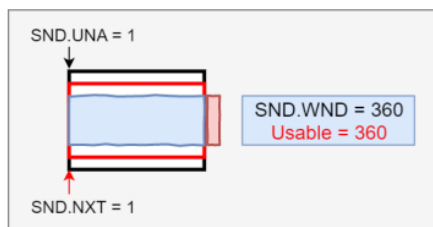
当应用程序没有及时读取缓存时，发送窗口和接收窗口的变化。

考虑以下场景：

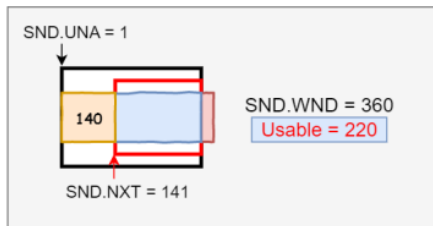
- 客户端作为发送方，服务端作为接收方，发送窗口和接收窗口初始大小为 360；
- 服务端非常的繁忙，当收到客户端的数据时，应用层不能及时读取数据。

客户端
(发送方)

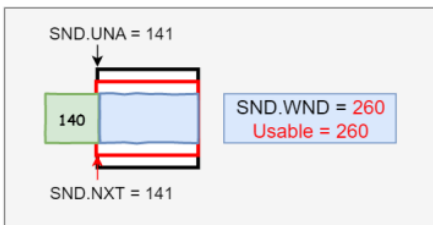
服务端
(接收方)



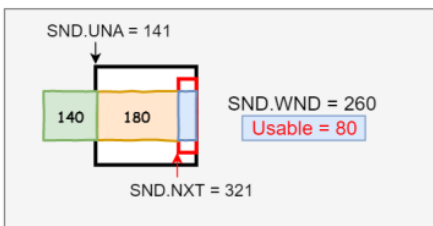
① 发送 140 字节的数据



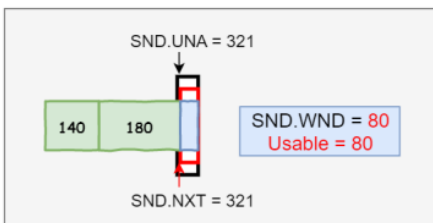
③ 收到接收方的窗口通告，发送窗口减少为 260



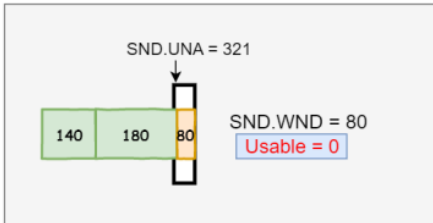
④ 发送 180 字节的数据



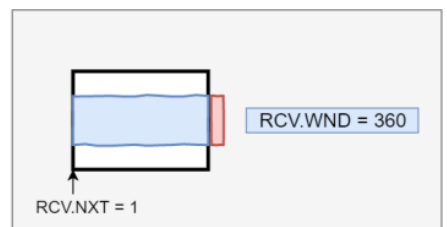
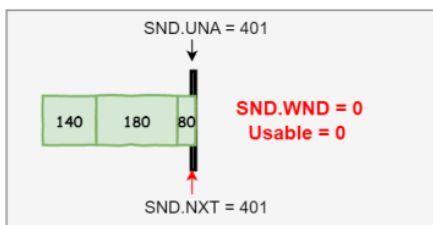
⑥ 收到接收方的窗口通告，发送窗口减少为 80



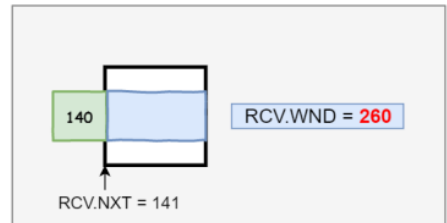
⑦ 发送 80 字节的数据



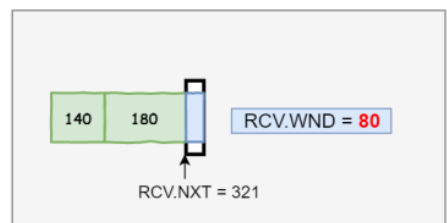
⑨ 收到接收方的窗口通告，发送窗口减少为 0



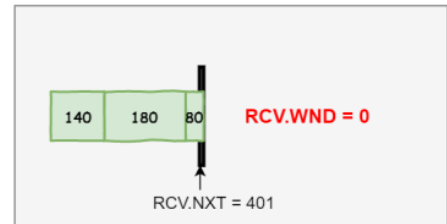
② 收到 140 字节数据后，但是应用进程只读取了 40 个字节，还有 100 字节一直占用着缓冲区，于是接收窗口收缩到了 260 (360 - 100)，并在发送确认信息时，通告窗口大小给发送方



⑤ 收到 180 字节数据后，但是应用进程读取任何数据，这 180 字节留在了缓冲区，于是接收窗口收缩到了 80 (260 - 180)，并在发送确认信息时，通告窗口大小给发送方



⑧ 收到 80 字节数据后，但是应用进程依然读取任何数据，这 80 字节留在了缓冲区，于是接收窗口收缩到了 0 (80 - 80)，并在发送确认信息时，通告窗口大小给发送方



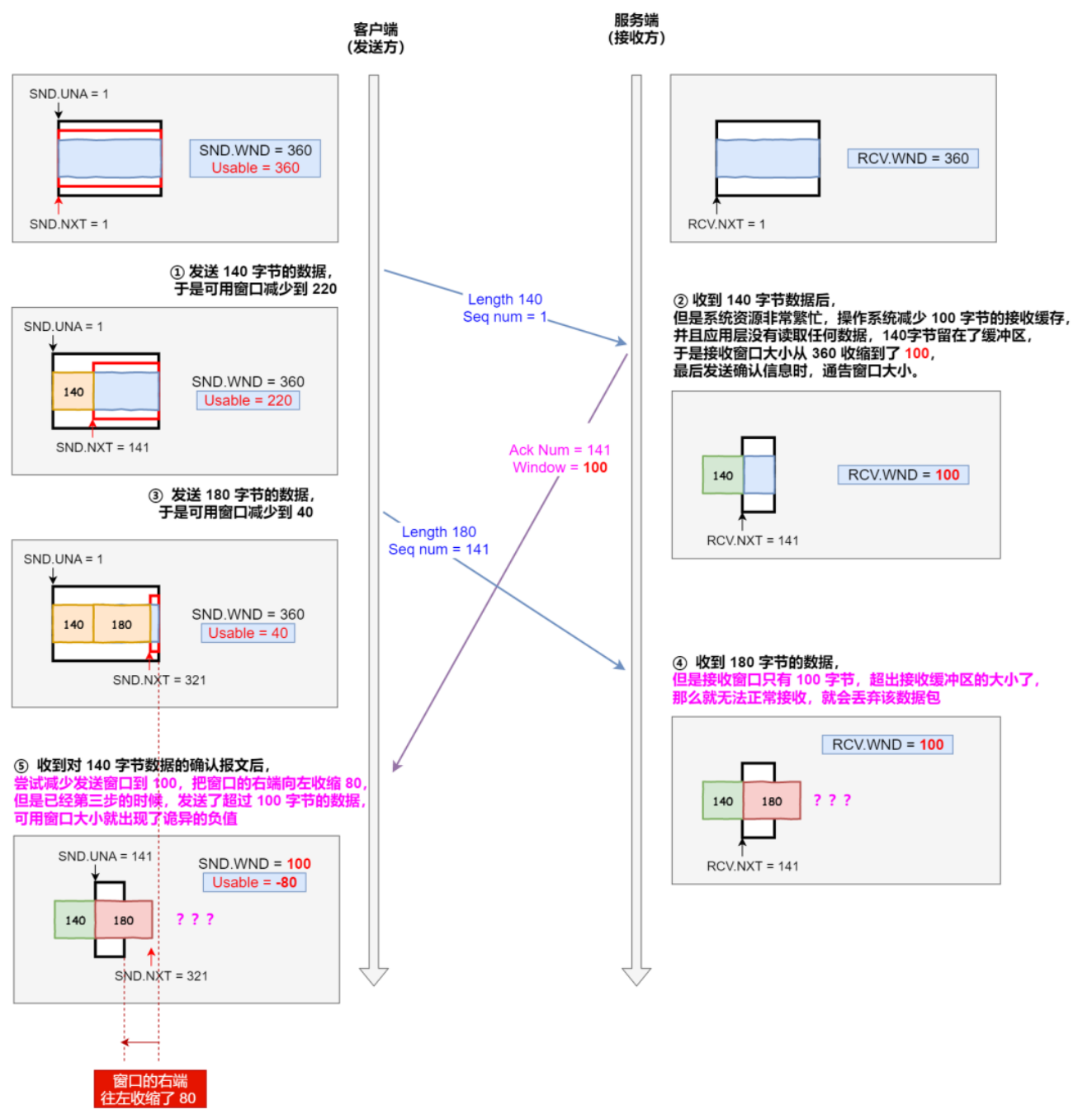
根据上图的流量控制，说明下每个过程：

1. 客户端发送 140 字节数据后，可用窗口变为 220 (360 - 140) 。
2. 服务端收到 140 字节数据，但是服务端非常繁忙，应用进程只读取了 40 个字节，还有 100 字节占用着缓冲区，于是接收窗口收缩到了 260 (360 - 100)，最后发送确认信息时，将窗口大小通过给客户端。
3. 客户端收到确认和窗口通告报文后，发送窗口减少为 260。
4. 客户端发送 180 字节数据，此时可用窗口减少到 80。
5. 服务端收到 180 字节数据，但是应用程序没有读取任何数据，这 180 字节直接就留在了缓冲区，于是接收窗口收缩到了 80 (260 - 180)，并在发送确认信息时，通过窗口大小给客户端。
6. 客户端收到确认和窗口通告报文后，发送窗口减少为 80。
7. 客户端发送 80 字节数据后，可用窗口耗尽。
8. 服务端收到 80 字节数据，但是应用程序依然没有读取任何数据，这 80 字节留在了缓冲区，于是接收窗口收缩到了 0，并在发送确认信息时，通过窗口大小给客户端。
9. 客户端收到确认和窗口通告报文后，发送窗口减少为 0。

可见最后窗口都收缩为 0 了，也就是发生了窗口关闭。当发送方可用窗口变为 0 时，发送方实际上会定时发送窗口探测报文，以便知道接收方的窗口是否发生了改变，这个内容后面会说，这里先简单提一下。

我们先来看看第二个例子。

当服务端系统资源非常紧张的时候，操心系统可能会直接减少了接收缓冲区大小，这时应用程序又无法及时读取缓存数据，那么这时候就有严重的事情发生了，会出现数据包丢失的现象。



说明下每个过程：

1. 客户端发送 140 字节的数据，于是可用窗口减少到了 220。
2. 服务端因为现在非常的繁忙，操作系统于是就把接收缓存减少了 100 字节，当收到对 140 数据确认报文后，又因为应用程序没有读取任何数据，所以 140 字节留在了缓冲区中，于是接收窗口大小从 360 收缩成了 100，最后发送确认信息时，通告窗口大小给对方。
3. 此时客户端因为还没有收到服务端的通告窗口报文，所以不知道此时接收窗口收缩成了 100，客户端只会看自己的可用窗口还有 220，所以客户端就发送了 180 字节数据，于是可用窗口减少到 40。
4. 服务端收到了 180 字节数据时，发现数据大小超过了接收窗口的大小，于是就把数据包丢失了。
5. 客户端收到第 2 步时，服务端发送的确认报文和通告窗口报文，尝试减少发送窗口到 100，把窗口的右端向左收缩了 80，此时可用窗口的大小就会出现诡异的负值。

所以，如果发生了先减少缓存，再收缩窗口，就会出现丢包的现象。

为了防止这种情况发生，TCP 规定是不允许同时减少缓存又收缩窗口的，而是采用先收缩窗口，过段时间在减少缓存，这样就可以避免了丢包情况。

窗口关闭

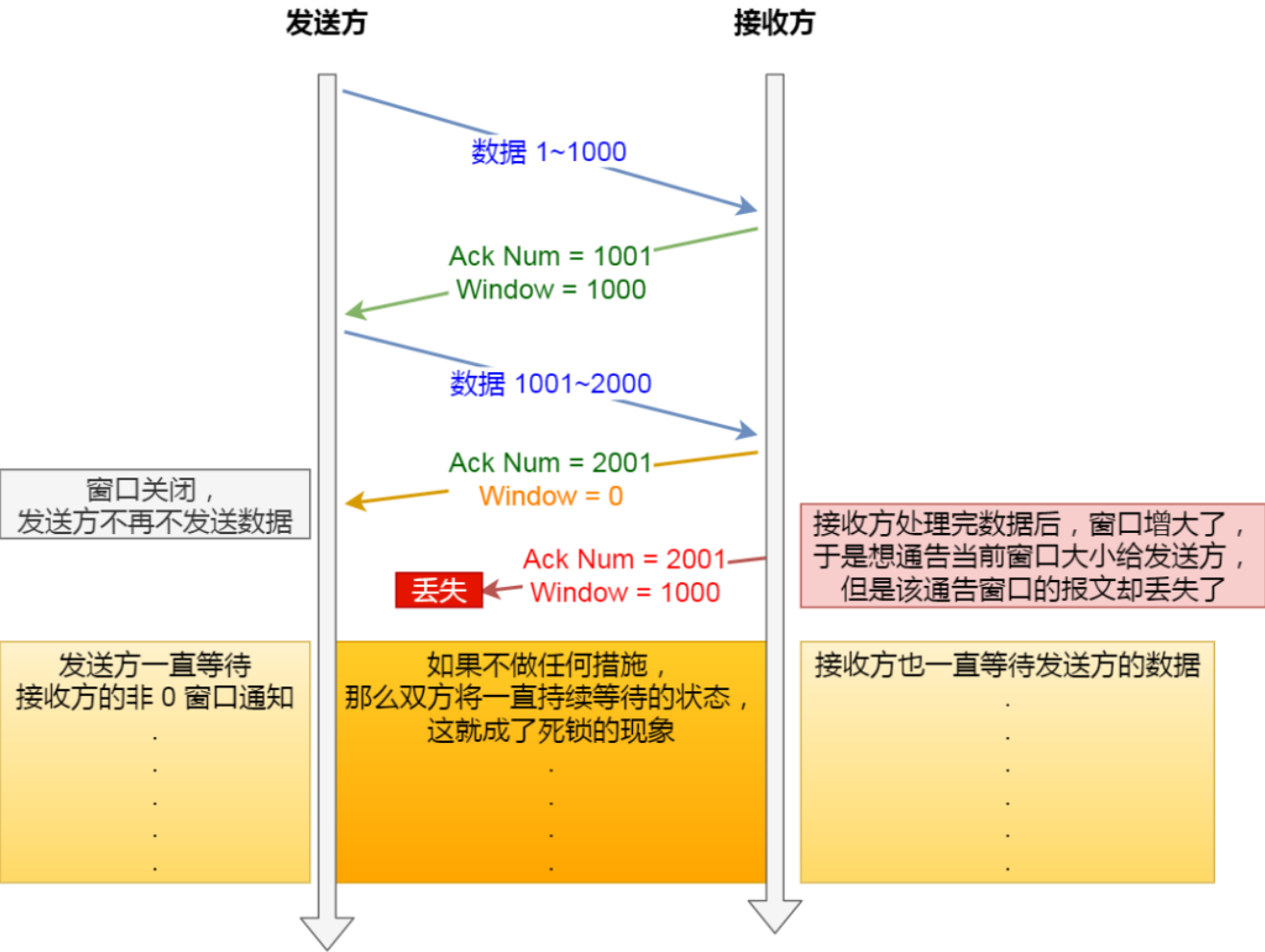
在前面我们都看到了，TCP 通过让接收方指明希望从发送方接收的数据大小（窗口大小）来进行流量控制。

如果窗口大小为 0 时，就会阻止发送方给接收方传递数据，直到窗口变为非 0 为止，这就是窗口关闭。

窗口关闭潜在的危险

接收方向发送方通告窗口大小时，是通过 ACK 报文来通告的。

那么，当发生窗口关闭时，接收方处理完数据后，会向发送方通告一个窗口非 0 的 ACK 报文，如果这个通告窗口的 ACK 报文在网络中丢失了，那麻烦就大了。

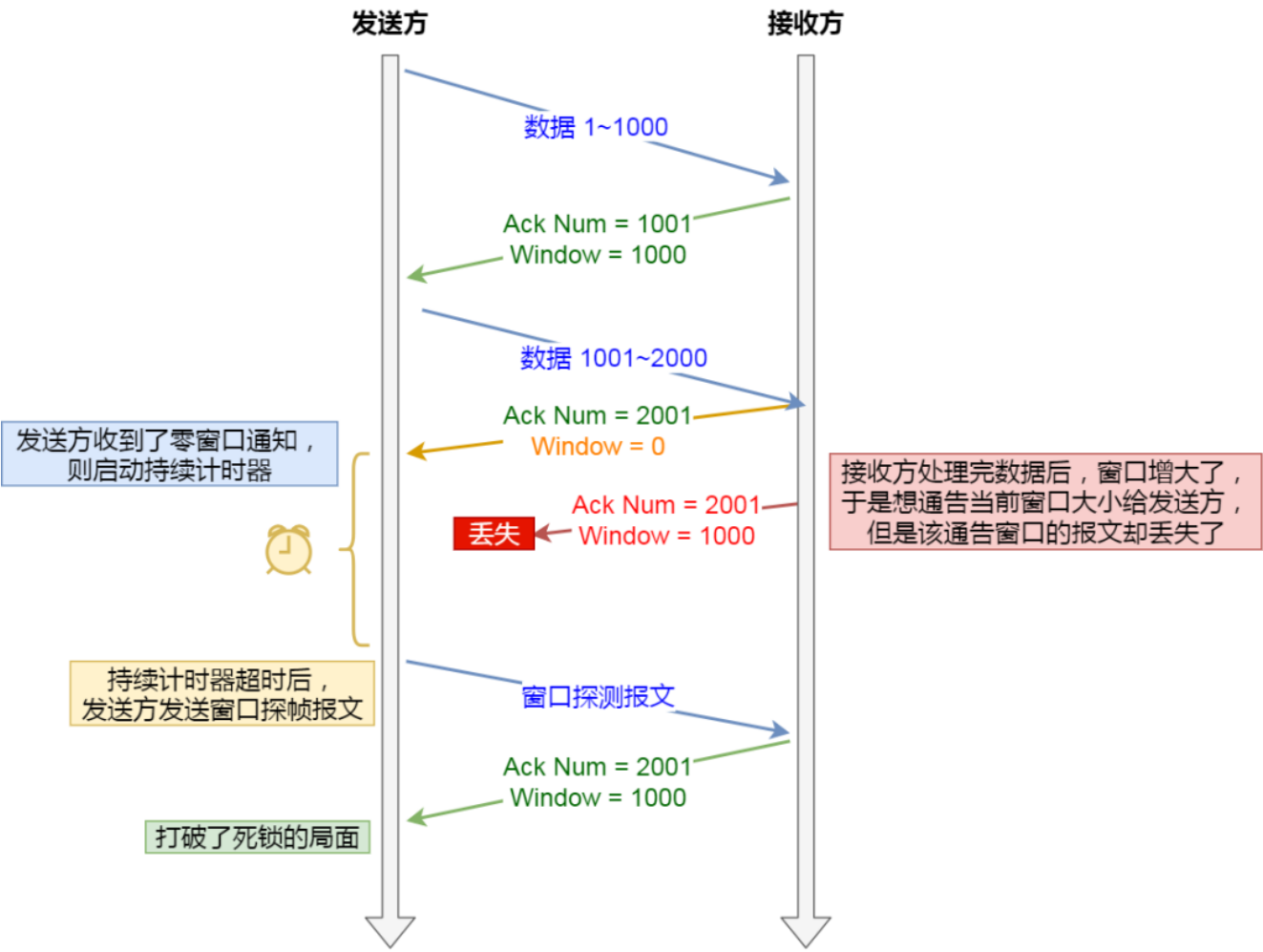


这会导致发送方一直等待接收方的非 0 窗口通知，接收方也一直等待发送方的数据，如不采取措施，这种相互等待的过程，会造成了死锁的现象。

TCP 是如何解决窗口关闭时，潜在的死锁现象呢？

为了解决这个问题，TCP 为每个连接设有一个持续定时器，只要 TCP 连接一方收到对方的零窗口通知，就启动持续计时器。

如果持续计时器超时，就会发送窗口探测 (Window probe) 报文，而对方在确认这个探测报文时，给出自己现在的接收窗口大小。



- 如果接收窗口仍然为 0，那么收到这个报文的一方就会重新启动持续计时器；
- 如果接收窗口不是 0，那么死锁的局面就可以被打破了。

窗口探查探测的次数一般为 3 次，每次大约 30-60 秒（不同的实现可能会不一样）。如果 3 次过后接收窗口还是 0 的话，有的 TCP 实现就会发 RST 报文来中断连接。

糊涂窗口综合症

如果接收方太忙了，来不及取走接收窗口里的数据，那么就会导致发送方的发送窗口越来越小。

到最后，如果接收方腾出几个字节并告诉发送方现在有几个字节的窗口，而发送方会义无反顾地发送这几个字节，这就是糊涂窗口综合症。

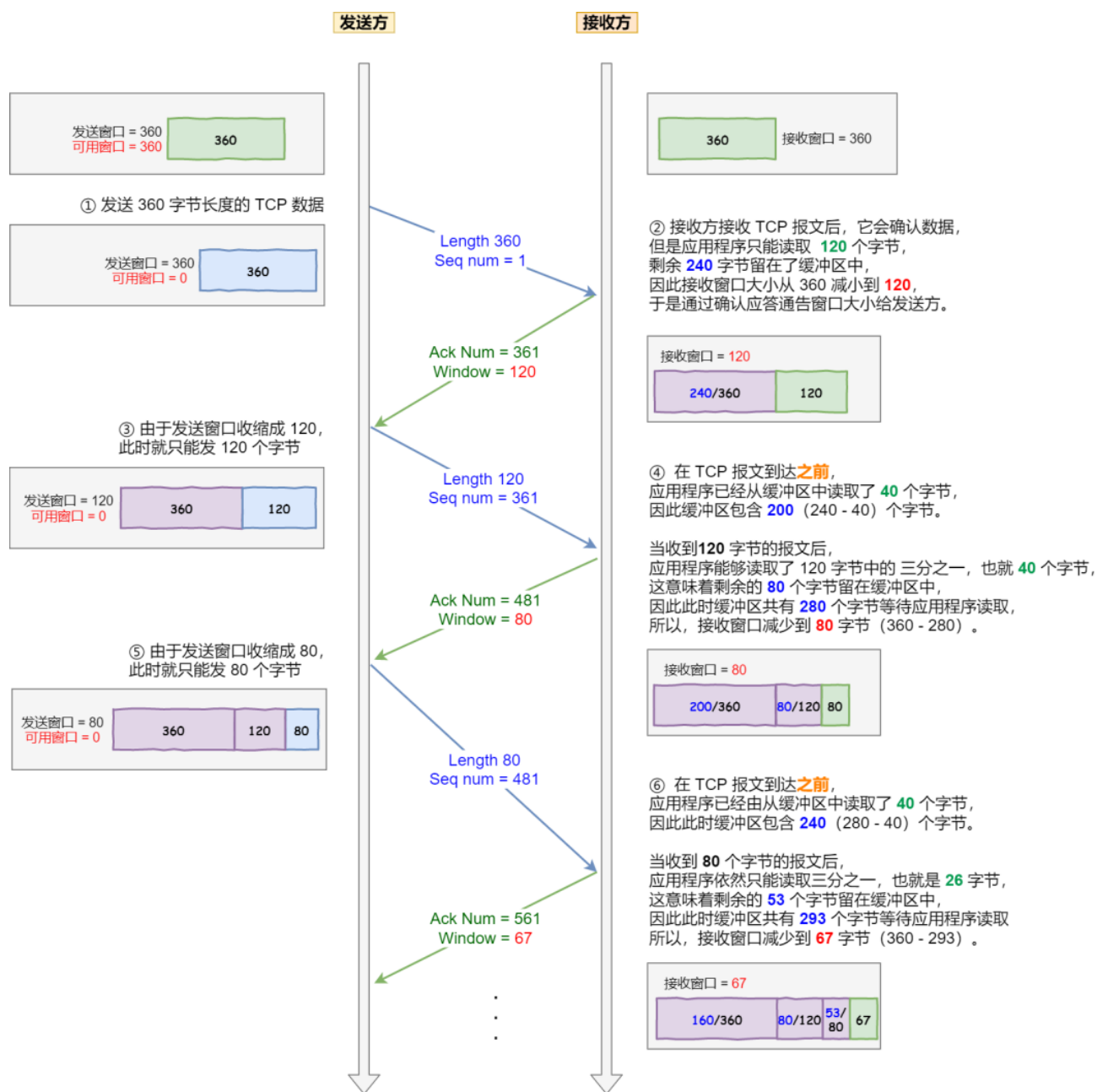
要知道，我们的 TCP + IP 头有 40 个字节，为了传输那几个字节的数据，要达上这么大的开销，这太不经济了。

就好像一个可以承载 50 人的大巴车，每次来了一两个人，就直接发车。除非家里有矿的大巴司机，才敢这样玩，不然迟早破产。要解决这个问题也不难，大巴司机等乘客数量超过了 25 个，才认定可以发车。

现举个糊涂窗口综合症的栗子，考虑以下场景：

接收方的窗口大小是 360 字节，但接收方由于某些原因陷入困境，假设接收方的应用层读取的能力如下：

- 接收方每接收 3 个字节，应用程序就只能从缓冲区中读取 1 个字节的数
- 在下一个发送方的 TCP 段到达之前，应用程序还从缓冲区中读取了 40 个额外的字节；



糊涂窗口综合症

每个过程的窗口大小的变化，在图中都描述的很清楚了，可以发现窗口不断减少了，并且发送的数据都是比较小的了。

所以，糊涂窗口综合症的现象是可以发生在发送方和接收

方：

- 接收方可以通告一个小的窗口
- 而发送方可以发送小数据

于是，要解决糊涂窗口综合症，就解决上面两个问题就可以了

- 让接收方不通告小窗口给发送方
- 让发送方避免发送小数据

怎么让接收方不通告小窗口呢？

接收方通常的策略如下：

当「窗口大小」小于 $\min(\text{MSS}, \text{缓存空间}/2)$ ，也就是小于 MSS 与 $1/2$ 缓存大小中的最小值时，就会向发送方通告窗口为 0，也就阻止了发送方再发数据过来。

等到接收方处理了一些数据后，窗口大小 $\geq \text{MSS}$ ，或者接收方缓存空间有一半可以使用，就可以把窗口打开让发送方发送数据过来。

怎么让发送方避免发送小数据呢？

发送方通常的策略：

使用 Nagle 算法，该算法的思路是延时处理，它满足以下两个条件中的一条才可以发送数据：

- 要等到窗口大小 $\geq \text{MSS}$ 或是 数据大小 $\geq \text{MSS}$

- 收到之前发送数据的 ack 回包

只要没满足上面条件中的一条，发送方一直在囤积数据，直到满足上面的发送条件。

另外，Nagle 算法默认是打开的，如果对于一些需要小数据包交互的场景的程序，比如，telnet 或 ssh 这样的交互性比较强的程序，则需要关闭 Nagle 算法。

可以在 Socket 设置 `TCP_NODELAY` 选项来关闭这个算法（关闭 Nagle 算法没有全局参数，需要根据每个应用自己的特点来关闭）

```
setsockopt(sock_fd, IPPROTO_TCP, TCP_NODELAY, (char *)&valu
```

拥塞控制

为什么要有拥塞控制呀，不是有流量控制了吗？

前面的流量控制是避免「发送方」的数据填满「接收方」的缓存，但是并不知道网络的中发生了什么。

一般来说，计算机网络都处在一个共享的环境。因此也有可能因为其他主机之间的通信使得网络拥堵。

在网络出现拥堵时，如果继续发送大量数据包，可能会导致数据包时延、丢失等，这时 TCP 就会重传数据，但是一重传就会导致网络的负担更重，于是会导致更大的延迟以及更多的丢包，这个情况就会进入恶性循环被不断地放大....

所以，TCP 不能忽略网络上发生的事，它被设计成一个无私的协议，当网络发送拥塞时，TCP 会自我牺牲，降低发送的数据量。

于是，就有了拥塞控制，控制的目的是避免「发送方」的数据填满整个网络。

为了在「发送方」调节所要发送数据的量，定义了一个叫做「拥塞窗口」的概念。

什么是拥塞窗口？和发送窗口有什么关系呢？

拥塞窗口 cwnd是发送方维护的一个的状态变量，它会根据网络的拥塞程度动态变化的。

我们在前面提到过发送窗口 `swnd` 和接收窗口 `rwnd` 是约等于的关系，那么由于入了拥塞窗口的概念后，此时发送窗口的值是 $swnd = \min(cwnd, rwnd)$ ，也就是拥塞窗口和接收窗口中的最小值。

拥塞窗口 `cwnd` 变化的规则：

- 只要网络中没有出现拥塞，`cwnd` 就会增大；
- 但网络中出现了拥塞，`cwnd` 就减少；

那么怎么知道当前网络是否出现了拥塞呢？

其实只要「发送方」没有在规定时间内接收到 ACK 应答报文，也就是发生了超时重传，就会认为网络出现了拥塞。

拥塞控制有哪些控制算法？

拥塞控制主要是四个算法：

- 慢启动
- 拥塞避免
- 拥塞发生
- 快速恢复

慢启动

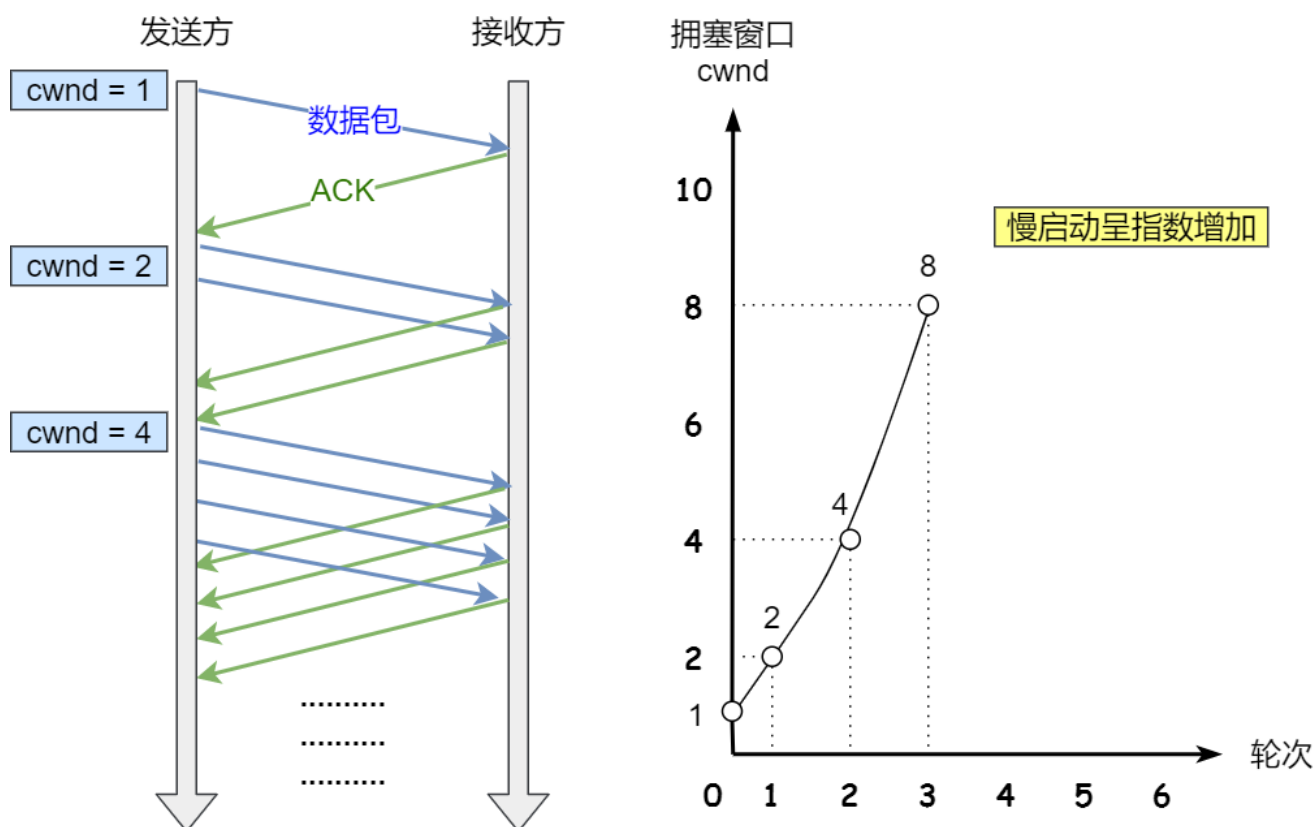
TCP 在刚建立连接完成后，首先是有个慢启动的过程，这个慢启动的意思就是一点一点的提高发送数据包的数量，如果一上来就发大量的数据，这不是给网络添堵吗？

慢启动的算法记住一个规则就行：当发送方每收到一个 **ACK**，就拥塞窗口 **cwnd** 的大小就会加 1。

这里假定拥塞窗口 **cwnd** 和发送窗口 **swnd** 相等，下面举个栗子：

- 连接建立完成后，一开始初始化 $cwnd = 1$ ，表示可以传一个 **MSS** 大小的数据。
- 当收到一个 **ACK** 确认应答后，**cwnd** 增加 1，于是一次能够发送 2 个
- 当收到 2 个的 **ACK** 确认应答后，**cwnd** 增加 2，于是就可以比之前多发 2 个，所以这一次能够发送 4 个

- 当这 4 个的 ACK 确认到来的时候，每个确认 cwnd 增加 1，4 个确认 cwnd 增加 4，于是就可以比之前多发 4 个，所以这一次能够发送 8 个。



慢启动算法

可以看出慢启动算法，发包的个数是指数性的增长。

那慢启动涨到什么时候是个头呢？

有一个叫慢启动门限 **ssthresh** (slow start threshold) 状态变量。

- 当 $\text{cwnd} < \text{ssthresh}$ 时，使用慢启动算法。
- 当 $\text{cwnd} \geq \text{ssthresh}$ 时，就会使用「拥塞避免算法」。

拥塞避免算法

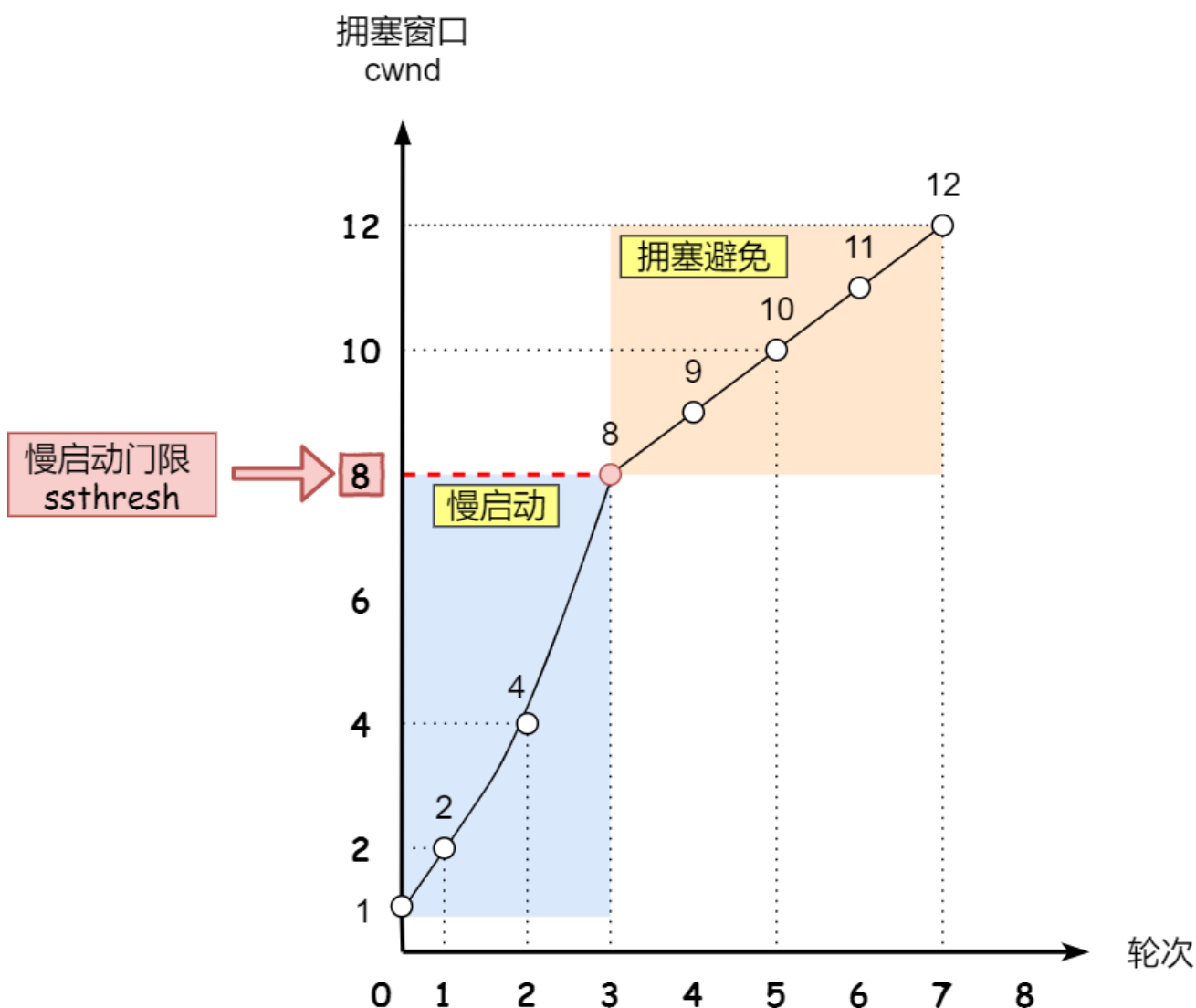
前面说道，当拥塞窗口 $cwnd$ 「超过」慢启动门限 $ssthresh$ 就会进入拥塞避免算法。

一般来说 $ssthresh$ 的大小是 65535 字节。

那么进入拥塞避免算法后，它的规则是：每当收到一个 **ACK** 时， $cwnd$ 增加 $1/cwnd$ 。

接上前面的慢启动的栗子，现假定 $ssthresh$ 为 8：

- 当 8 个 ACK 应答确认到来时，每个确认增加 $1/8$ ，8 个 ACK 确认 $cwnd$ 一共增加 1，于是这一次能够发送 9 个 MSS 大小的数据，变成了线性增长。



所以，我们可以发现，拥塞避免算法就是将原本慢启动算法的指数增长变成了线性增长，还是增长阶段，但是增长速度缓慢了一些。

就这么一直增长着后，网络就会慢慢进入了拥塞的状况了，于是就会出现丢包现象，这时就需要对丢失的数据包进行重传。

当触发了重传机制，也就进入了「拥塞发生算法」。

拥塞发生

当网络出现拥塞，也就是会发生数据包重传，重传机制主要有两种：

- 超时重传
- 快速重传

这两种使用的拥塞发送算法是不同的，接下来分别来说。

发生超时重传的拥塞发生算法

当发生了「超时重传」，则就会使用拥塞发生算法。

这个时候，`ssthresh` 和 `cwnd` 的值会发生变化：

- `ssthresh` 设为 $cwnd/2$,
- `cwnd` 重置为 1

接着，就重新开始慢启动，慢启动是会突然减少数据流的。这真是一旦「超时重传」，马上回到解放前。但是这种方式太激进了，反应也很强烈，会造成网络卡顿。

就好像本来在秋名山高速漂移着，突然来个紧急刹车，轮胎受得了吗。。。

发生快速重传的拥塞发生算法

还有更好的方式，前面我们讲过「快速重传算法」。当接收方发现丢了一个中间包的时候，发送三次前一个包的ACK，于是发送端就会快速地重传，不必等待超时再重传。

TCP 认为这种情况不严重，因为大部分没丢，只丢了一小部分，则 `ssthresh` 和 `cwnd` 变化如下：

- $cwnd = cwnd / 2$ ，也就是设置为原来的一半；
- `ssthresh = cwnd`;
- 进入快速恢复算法

快速恢复

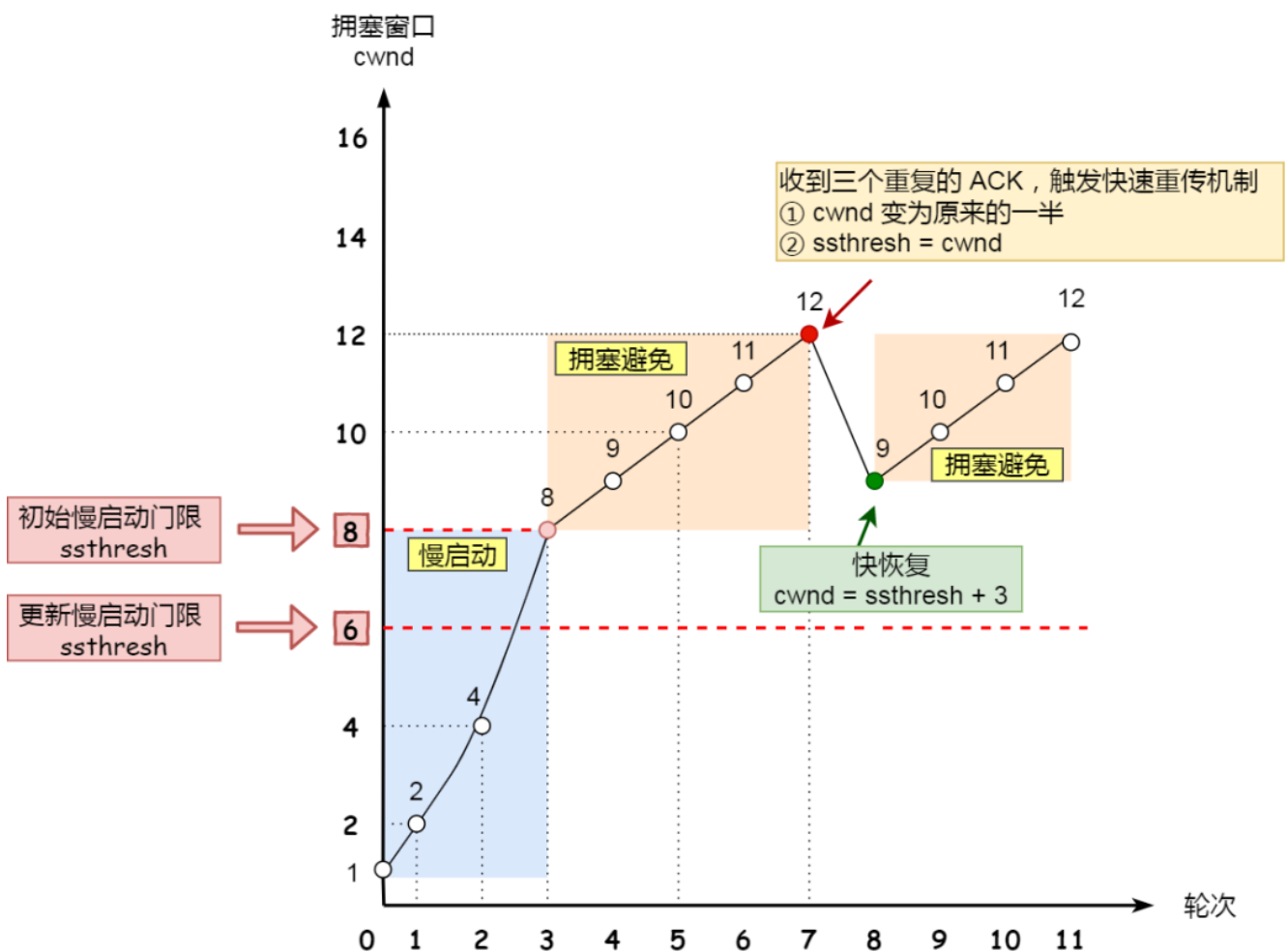
快速重传和快速恢复算法一般同时使用，快速恢复算法是认为，你还能收到 3 个重复 ACK 说明网络也不那么糟糕，所以没有必要像 `RTO` 超时那么强烈。

正如前面所说，进入快速恢复之前，`cwnd` 和 `ssthresh` 已被更新了：

- $\text{cwnd} = \text{cwnd}/2$, 也就是设置为原来的一半;
- $\text{ssthresh} = \text{cwnd}$;

然后, 进入快速恢复算法如下:

- 拥塞窗口 $\text{cwnd} = \text{ssthresh} + 3$ (3 的意思是确认有 3 个数据包被收到了)
- 重传丢失的数据包
- 如果再收到重复的 ACK, 那么 cwnd 增加 1
- 如果收到新数据的 ACK 后, 设置 cwnd 为 ssthresh , 接着就进入了拥塞避免算法



也就是没有像「超时重传」一夜回到解放前，而是还在比较高的值，后续呈线性增长。

巨人的肩膀

[1] 趣谈网络协议专栏.刘超.极客时间

[2] Web协议详解与抓包实战专栏.陶辉.极客时间

[3] TCP/IP详解 卷1：协议.范建华 译.机械工业出版社

[4] 图解TCP/IP.竹下隆史.人民邮电出版社

[5] The TCP/IP Guide.Charles M. Kozierok.

[6] TCP那些事（上）.陈皓.酷壳博客.
<https://coolshell.cn/articles/11564.html>

[7] TCP那些事（下）.陈皓.酷壳博客.
<https://coolshell.cn/articles/11609.html>

往期好文

[硬不硬你说了算！近 40 张图解被问千百遍的 TCP 三次握手和四次挥手面试题](#)

[硬核！30 张图解 HTTP 常见的面试题](#)

唠叨唠叨

是吧？TCP 巨复杂吧？看完很累吧？

但这还只是 TCP 冰山一脚，它的更深处就由你们自己去探索啦。

送给读者大大们一个回眸一笑，很倾城的吧

本文只是抛砖引玉，若你有更好的想法或文章有误的地方，欢迎留言讨论！

小林是专为大家图解的工具人，**Goodbye**，我们下次见！