

TCP:

Control de Congestión

Introducción a los Sistemas Distribuidos (75.43)

Universidad de Buenos Aires, Facultad de Ingeniería



Control de flujo

Control de flujo

- Cada extremo de la transmisión TCP tiene su propio buffer.
- TCP almacena datos en el buffer para pasarlos en orden a la capa de aplicación.
- El buffer tiene un tamaño finito.



Control de flujo

¿Qué pasa si se reciben más datos de los que la capa de aplicación es capaz de consumir?

Control de flujo

¿Qué pasa si se reciben más datos de los que la capa de aplicación es capaz de consumir?

BUFFER OVERFLOW

Control de flujo

- Para evitar el buffer overflow, **cada extremo** de la transmisión tiene que comunicar la cantidad de datos que puede recibir.
- A esto se lo llama ***rwnd*** (Receiver Window o Advertised Window).
- Se utiliza el campo Window del TCP header para comunicar este valor.

Control de flujo

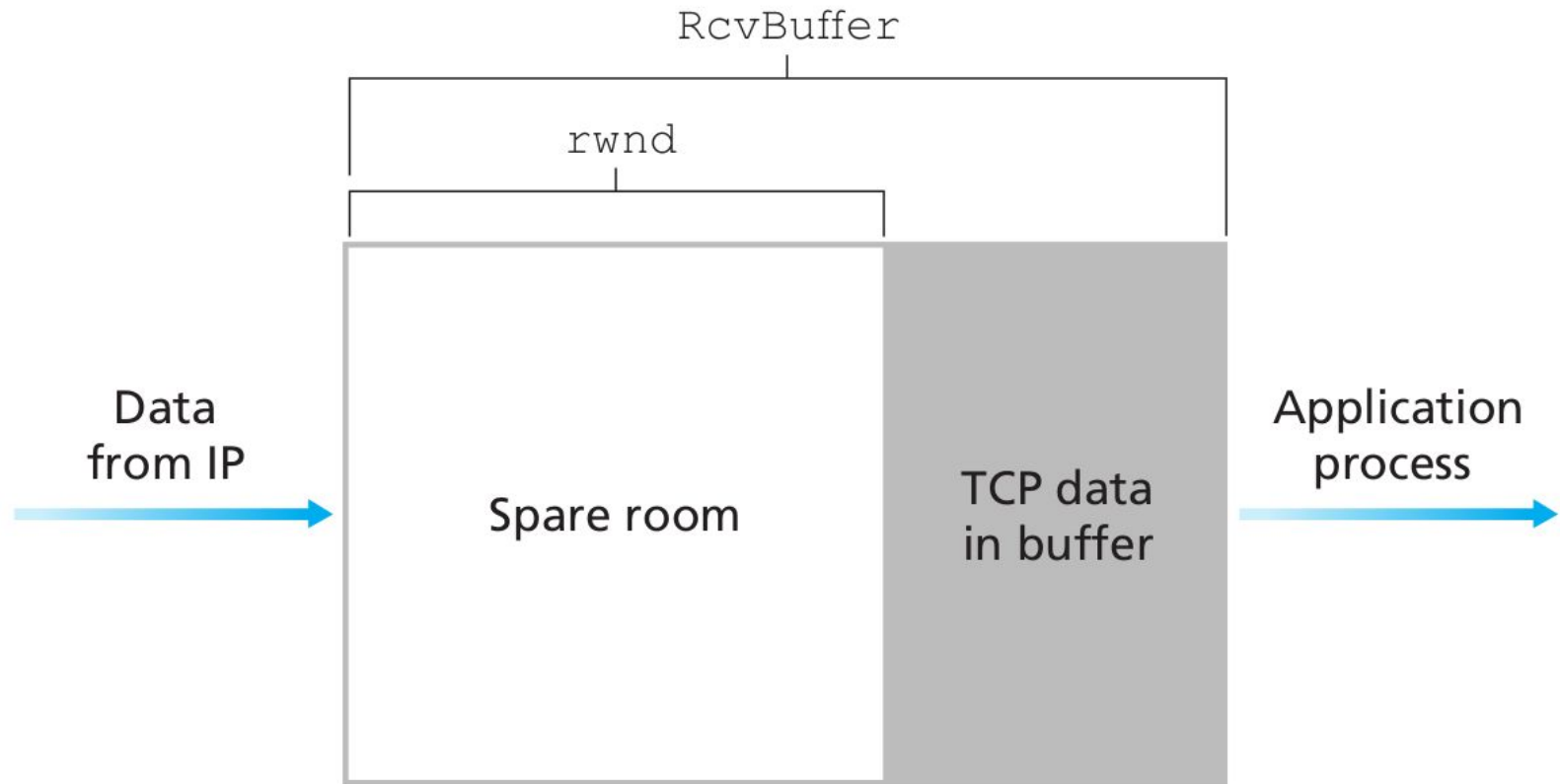


Figure 3.38 ♦ The receive window (**rwnd**) and the receive buffer (**RcvBuffer**)

Control de flujo

- Cuando se recibe un cambio en ***rwnd***, el emisor debe modificar su tasa de envío para no causar buffer overflow, pero tampoco subutilizar el buffer.
- Si se recibe ***rwnd*** = 0, el emisor debe dejar de enviar datos hasta recibir un nuevo valor de ***rwnd***

Control de flujo

- Imaginemos el siguiente escenario
 - A envía un paquete a B
 - El buffer de B se llena, por lo que envía un ACK *rwnd* = 0 a A.
 - A deja de enviar datos.



Control de flujo

- La capa de aplicación en B lee del buffer y lo libera.

¿Cómo hace B para avisarle a A que puede recibir más datos?



Control de flujo

- Para lograr esto, luego de un período de inactividad A continuará enviando paquetes de 1 byte a B.
- Cuando B los recibe, indica el tamaño de ventana actual.
- Cuando el buffer se libere, B indicará ***rwnd*** > 0 y se reanudará la transmisión.



Control de congestión

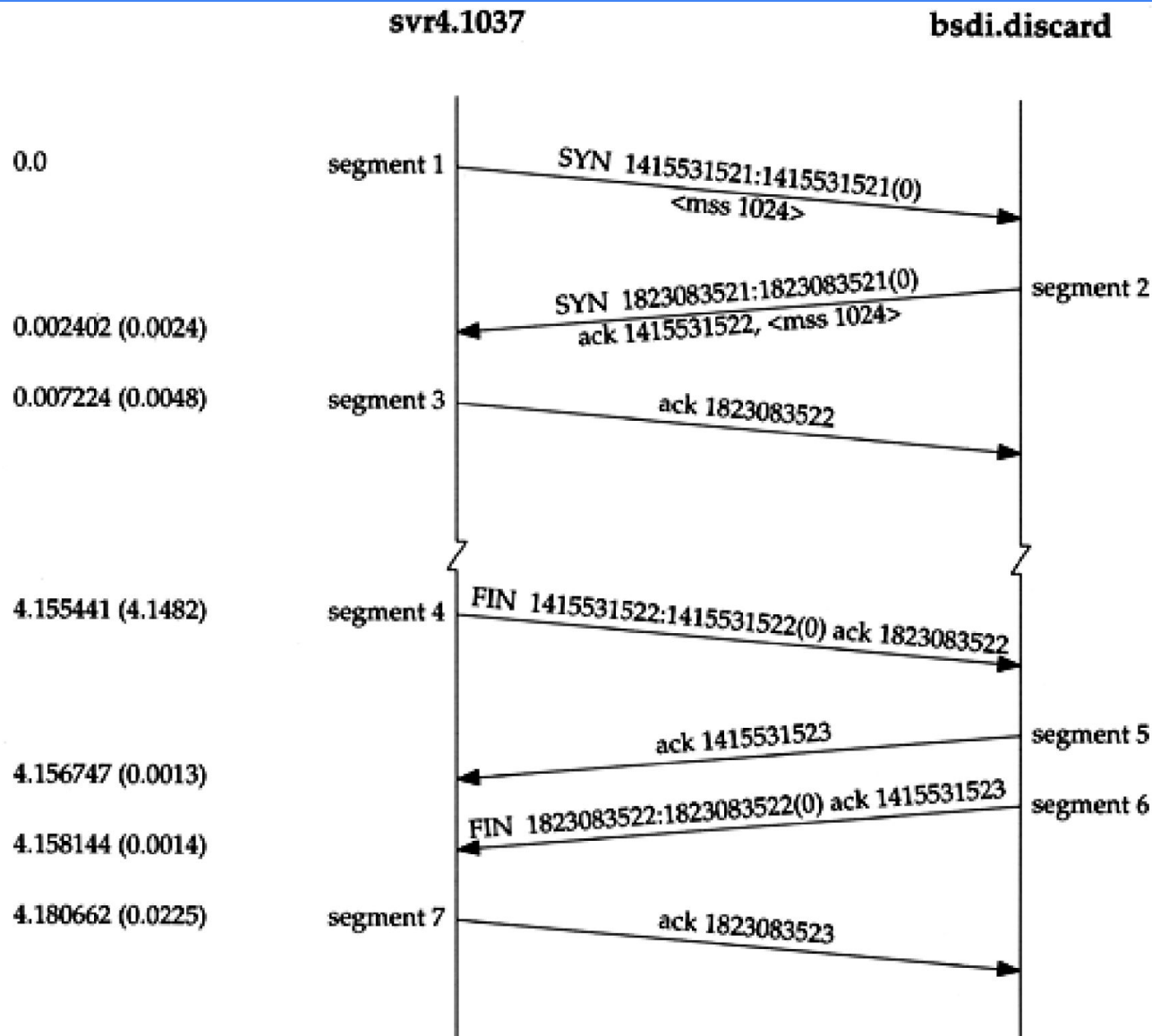
TCP : control de congestión

- El control de congestión **no formaba parte** del diseño original de TCP

TCP : control de congestión

- Fue una respuesta a la **congestión persistente** observada en la red
- La **congestión** es producto de múltiples flujos TCP en simultáneo sobre la misma red

TCP : control de congestión



Algoritmos de control de congestión



TCP : control de congestión

- La congestión aparece principalmente a causa del **tamaño finito** de los buffers de los **dispositivos intermedios** de capas inferiores.
- Cuando se llena un buffer, se empiezan a **descartar** los segmentos
-> se dice que la red está **congestionada**

TCP : control de congestión

- **cwnd** : es la máxima cantidad de bytes que puede haber en vuelo

- **RECORDAR:**

$$\text{LastByteSent} - \text{LastByteACKed} < \min(\text{cwnd}, \text{rwnd})$$

TCP : control de congestión

- Tiene distintas etapas:
 - Slow start
 - Congestion avoidance
 - Fast Retransmit
 - Fast Recovery

TCP : control de congestión

Slow start

- $\text{cwnd}(n+1) = \text{cwnd}(n) + \text{MSS} * \#(\text{ACK})$
- cwnd medido en Bytes

TCP : control de congestión

Slow start

- $\text{cwnd}(n+1) = \text{cwnd}(n) + \#(\text{ACK})$
- cwnd medido en **MSSs** (segmentos de tamaño máximo)

TCP : control de congestión

Slow start

- $\text{cwnd}(n+1) = \text{cwnd}(n) + \#(\text{ACK})$
- cwnd medido en **MSSs** (segmentos de tamaño máximo)
- Crecimiento exponencial

TCP : control de congestión

Slow start

$$cwnd(n+1) = cwnd(n) + \#(ACK)$$

- Ejemplo: si $cwnd(0) = 1$
- $cwnd(0) = 1$, envió 1 segmento

TCP : control de congestión

Slow start

$$cwnd(n+1) = cwnd(n) + \#(ACK)$$

- Ejemplo: si $cwnd(0) = 1$
- $cwnd(0) = 1$, envió 1 segmento
- Recibo 1 ACK
- $cwnd(1) = ?$

TCP : control de congestión

Slow start

$$cwnd(n+1) = cwnd(n) + \#(ACK)$$

- Ejemplo: si $cwnd(0) = 1$
- $cwnd(0) = 1$, envió 1 segmento
- Recibo 1 ACK
- $cwnd(1) = 1 + \#(ACK) = 1 + 1 = 2$

TCP : control de congestión

Slow start

$$cwnd(n+1) = cwnd(n) + \#(ACK)$$

- Ejemplo: si $cwnd(0) = 1$
- $cwnd(0) = 1$, envío 1 segmento
- Recibo 1 ACK
- $cwnd(1) = 1 + \#(ACK) = 1 + 1 = 2$
- Envío 2 segmentos. Recibo 2 ACK
- $cwnd(2) = ?$

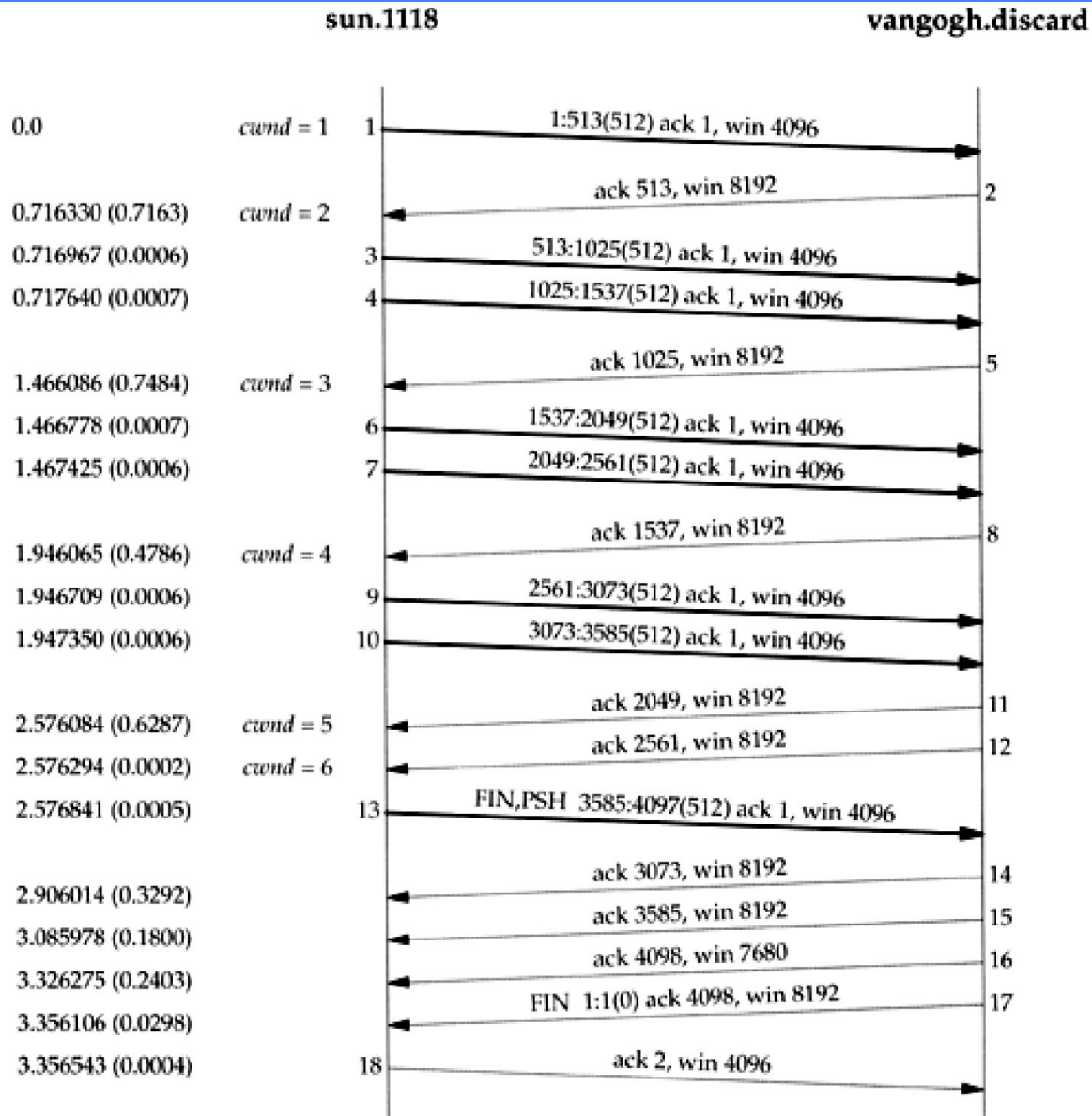
TCP : control de congestión

Slow start

$$cwnd(n+1) = cwnd(n) + \#(ACK)$$

- Ejemplo: si $cwnd(0) = 1$
- $cwnd(0) = 1$, envío 1 segmento
- Recibo 1 ACK
- $cwnd(1) = 1 + \#(ACK) = 1 + 1 = 2$
- Envío 2 segmentos. Recibo 2 ACK
- $cwnd(2) = 2 + \#(ACK) = 2 + 2 = 4$

TCP : control de congestión



TCP : control de congestión

¿Seguimos
incrementando la ventana
exponencialmente para
siempre?

TCP : control de congestión

Cambio de etapa

- En un determinado momento, se dará por terminada la etapa de **Slow Start** para pasar a la próxima, llamada **Congestion Avoidance**
- Esto viene dado por el valor de **ssthresh** (slow start threshold size), qué es configurable.

TCP : control de congestión

Congestion avoidance

- $\text{cwnd}(n+1) = \text{cwnd}(n) + \#(\text{ACK})/\text{cwnd}(n)$
- Aumenta 1 cuando nos llegan los ACKs de toda la ráfaga
- cwnd siempre tiene que ser entero (la unidad es el **MSS**)
- Se redondea para abajo!

TCP : control de congestión

Congestion avoidance

$$cwnd(n+1) = cwnd(n) + \#(ACK)/cwnd(n)$$

- Ejemplo: si $cwnd(n) = 2$
- $cwnd(n) = 2$, envió 2 segmentos

TCP : control de congestión

Congestion avoidance

$$cwnd(n+1) = cwnd(n) + \#(ACK)/cwnd(n)$$

- Ejemplo: si $cwnd(n) = 2$
- $cwnd(n) = 2$, envió 2 segmentos
- Llega 1 ACK
- $cwnd(n+1) = 2 + 1 / 2 = 2$

TCP : control de congestión

Congestion avoidance

$$cwnd(n+1) = cwnd(n) + \#(ACK)/cwnd(n)$$

- Ejemplo: si $cwnd(n) = 2$
- $cwnd(n) = 2$, envió 2 segmentos
- Llega 1 ACK
- $cwnd(n+1) = 2 + 1 / 2 = 2$
- Llega el 2do ACK
- $cwnd(n+2) = 2 + 2 / 2 = 3$

TCP : control de congestión

¿Seguimos en CA para
siempre, felizmente
incrementando en
 $\#(\text{ACK})/\text{cwnd}(n)$ por toda
la eternidad?

TCP : control de congestión

¿Qué pasa si se pierden
paquetes?

TCP : control de congestión

Pérdida por timeout (RTO)

- $ssthresh = cwnd(n) / 2$
- $cwnd(n+1) = 1$ (1 es LW [loss window])
- Empieza slow start de nuevo
- Ver ejemplo

TCP : control de congestión

En una conexión TCP recién establecida ($IW = 2$ MSS, $SSTHRESH = 64$ KB) con $RTT=200$ ms y $MSS=2$ KB, el host receptor siempre anuncia una AdvertisedWindow de 16 KB. La red está cargada al punto que si una ráfaga fuera de 16 KB o más, se perderían todos los segmentos de la misma. **¿Cuántos rounds se demora en enviar un archivo de 36 KB?**

TCP : control de congestión

En una conexión TCP recién establecida ($IW = 2$ MSS, $SSTHRESH = 64$ KB) con $RTT=200$ ms y $MSS=2$ KB, el host receptor siempre anuncia una AdvertisedWindow de 16 KB. La red está cargada al punto que si una ráfaga fuera de 16 KB o más, se perderían todos los segmentos de la misma. **¿Cuántos rounds se demora en enviar un archivo de 36 KB?**

¿Qué datos son importantes?

TCP : control de congestión

En una conexión TCP recién establecida ($IW = 2$ MSS, $SSTHRESH = 64$ KB) con $RTT=200$ ms y $MSS=2$ KB el host receptor siempre anuncia una AdvertisedWindow de 16 KB. La red está cargada al punto que si una ráfaga fuera de 16 KB o más, se perderían todos los segmentos de la misma. ¿Cuántos rounds se demora en enviar un archivo de 36 KB?

TCP : control de congestión

En una conexión TCP recién establecida ($IW = 2 \text{ MSS}$, $SSTHRESH = 32 \text{ MSS}$) con $RTT=200 \text{ ms}$ y $MSS=2 \text{ KB}$, el host receptor siempre anuncia una AdvertisedWindow de 8 MSS . La red está cargada al punto que si una ráfaga fuera de 8 MSS o más, se perderían todos los segmentos de la misma. ¿Cuántos rounds se demora en enviar un archivo de 18 MSS ?

Convertimos los datos a MSSs

Resolución

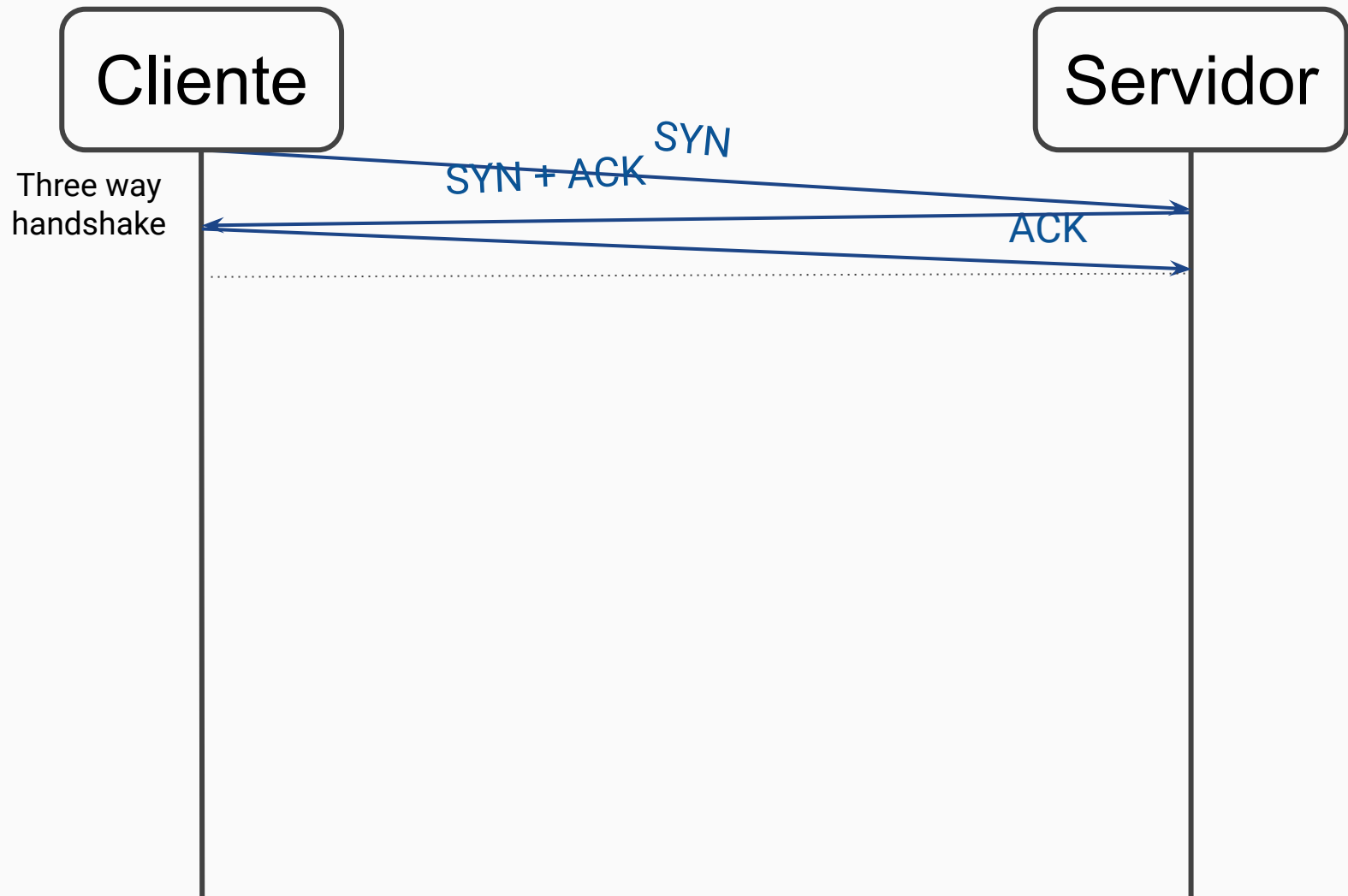
Cliente



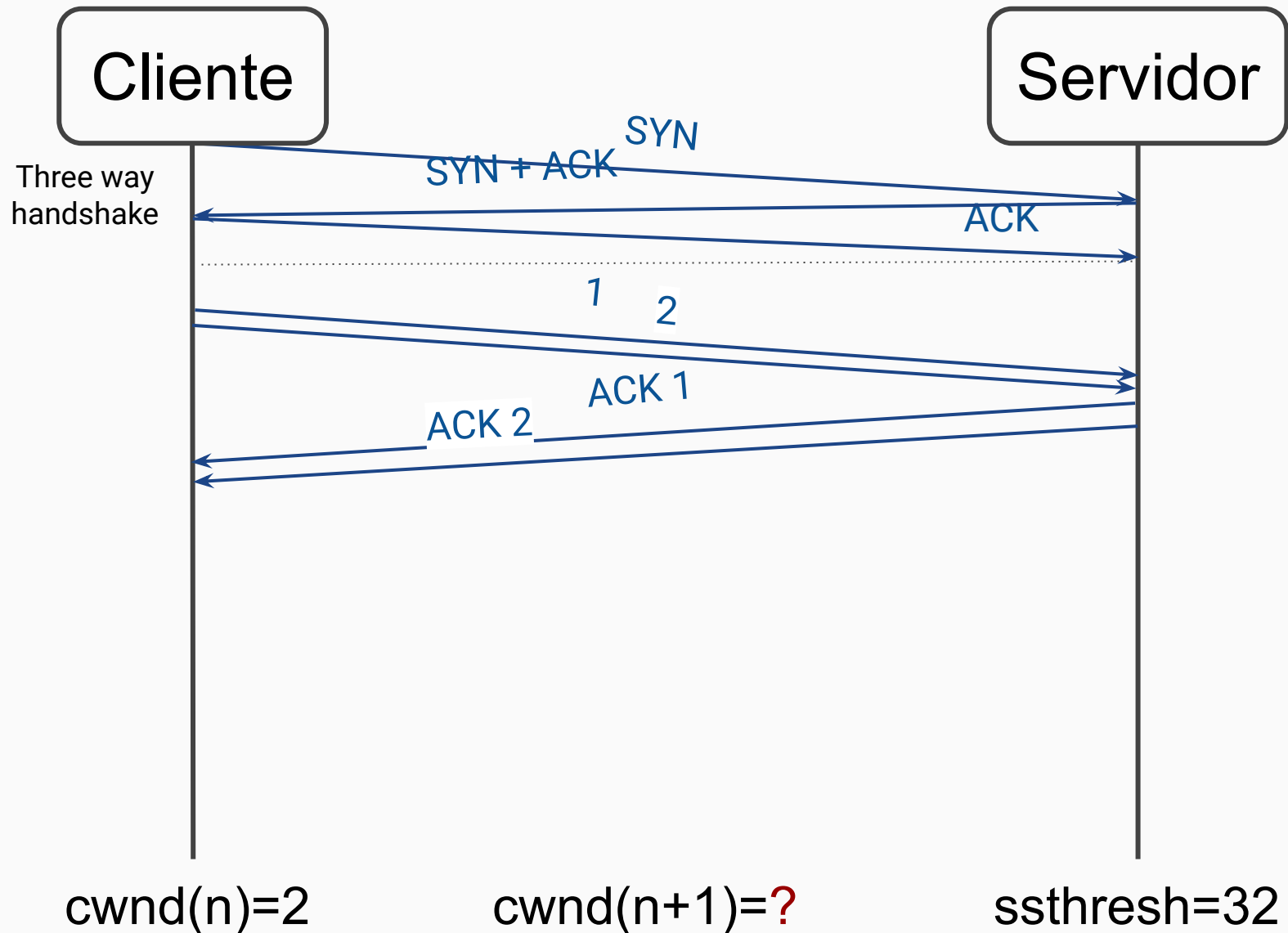
```
graph TD; C[Cliente] --- L1[ ]; L1 --- L2[ ]; L2 --- L3[ ]; L3 --- L4[ ]; L4 --- L5[ ]; L5 --- L6[ ]; L6 --- L7[ ]; L7 --- L8[ ]; L8 --- L9[ ]; L9 --- L10[ ]; L10 --- L11[ ]; L11 --- L12[ ]; L12 --- L13[ ]; L13 --- L14[ ]; L14 --- L15[ ]; L15 --- L16[ ]; L16 --- L17[ ]; L17 --- L18[ ]; L18 --- L19[ ]; L19 --- L20[ ]; L20 --- L21[ ]; L21 --- L22[ ]; L22 --- L23[ ]; L23 --- L24[ ]; L24 --- L25[ ]; L25 --- L26[ ]; L26 --- L27[ ]; L27 --- L28[ ]; L28 --- L29[ ]; L29 --- L30[ ]; L30 --- L31[ ]; L31 --- L32[ ]; L32 --- L33[ ]; L33 --- L34[ ]; L34 --- L35[ ]; L35 --- L36[ ]; L36 --- L37[ ]; L37 --- L38[ ]; L38 --- L39[ ]; L39 --- L40[ ]; L40 --- L41[ ]; L41 --- L42[ ]; L42 --- L43[ ]; L43 --- L44[ ]; L44 --- L45[ ]; L45 --- L46[ ]; L46 --- L47[ ]; L47 --- L48[ ]; L48 --- L49[ ]; L49 --- L50[ ]; L50 --- L51[ ]; L51 --- L52[ ]; L52 --- L53[ ]; L53 --- L54[ ]; L54 --- L55[ ]; L55 --- L56[ ]; L56 --- L57[ ]; L57 --- L58[ ]; L58 --- L59[ ]; L59 --- L60[ ]; L60 --- L61[ ]; L61 --- L62[ ]; L62 --- L63[ ]; L63 --- L64[ ]; L64 --- L65[ ]; L65 --- L66[ ]; L66 --- L67[ ]; L67 --- L68[ ]; L68 --- L69[ ]; L69 --- L70[ ]; L70 --- L71[ ]; L71 --- L72[ ]; L72 --- L73[ ]; L73 --- L74[ ]; L74 --- L75[ ]; L75 --- L76[ ]; L76 --- L77[ ]; L77 --- L78[ ]; L78 --- L79[ ]; L79 --- L80[ ]; L80 --- L81[ ]; L81 --- L82[ ]; L82 --- L83[ ]; L83 --- L84[ ]; L84 --- L85[ ]; L85 --- L86[ ]; L86 --- L87[ ]; L87 --- L88[ ]; L88 --- L89[ ]; L89 --- L90[ ]; L90 --- L91[ ]; L91 --- L92[ ]; L92 --- L93[ ]; L93 --- L94[ ]; L94 --- L95[ ]; L95 --- L96[ ]; L96 --- L97[ ]; L97 --- L98[ ]; L98 --- L99[ ]; L99 --- L100[ ]; L100 --- S[Servidor];
```

Servidor

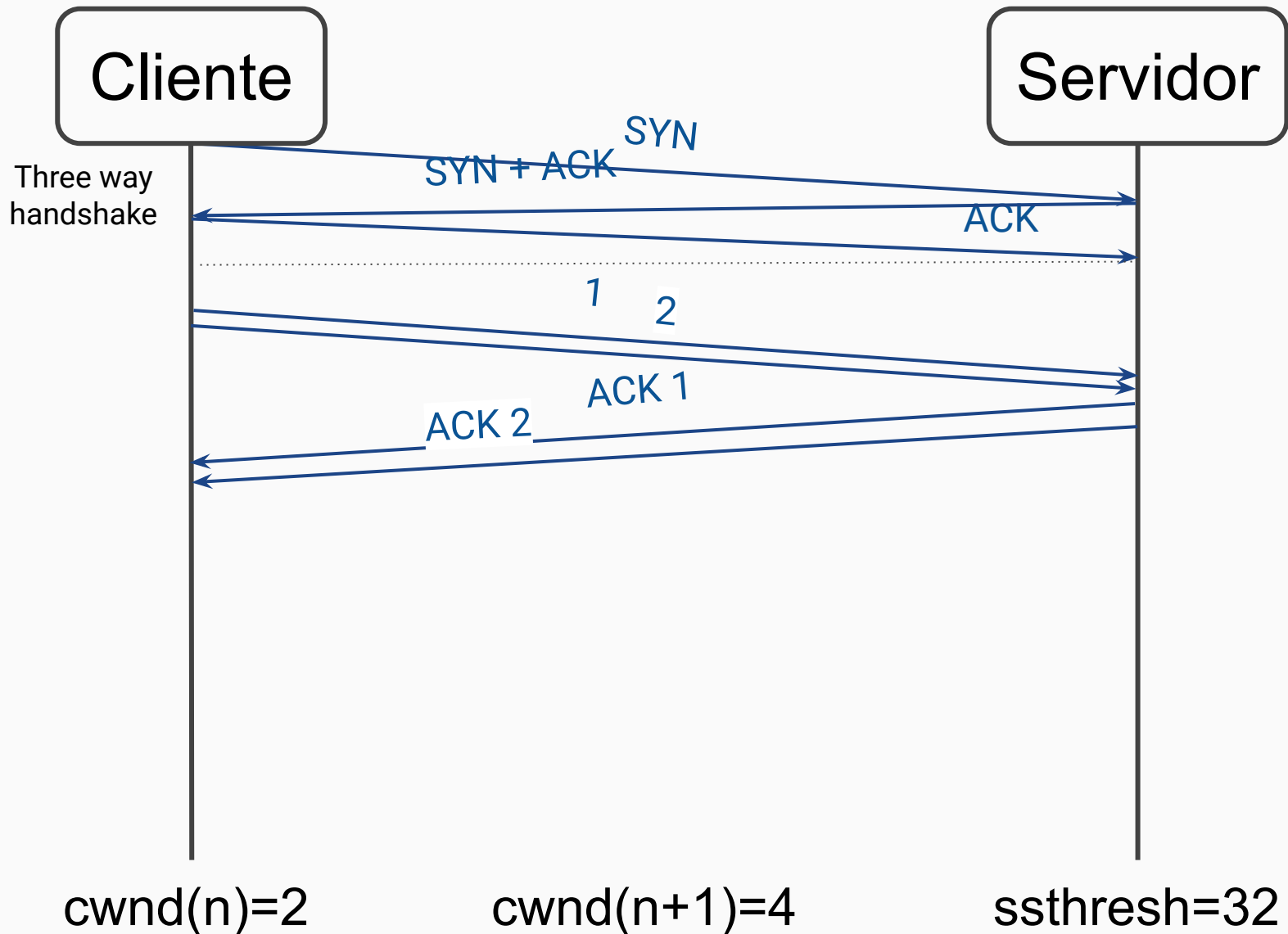
Resolución



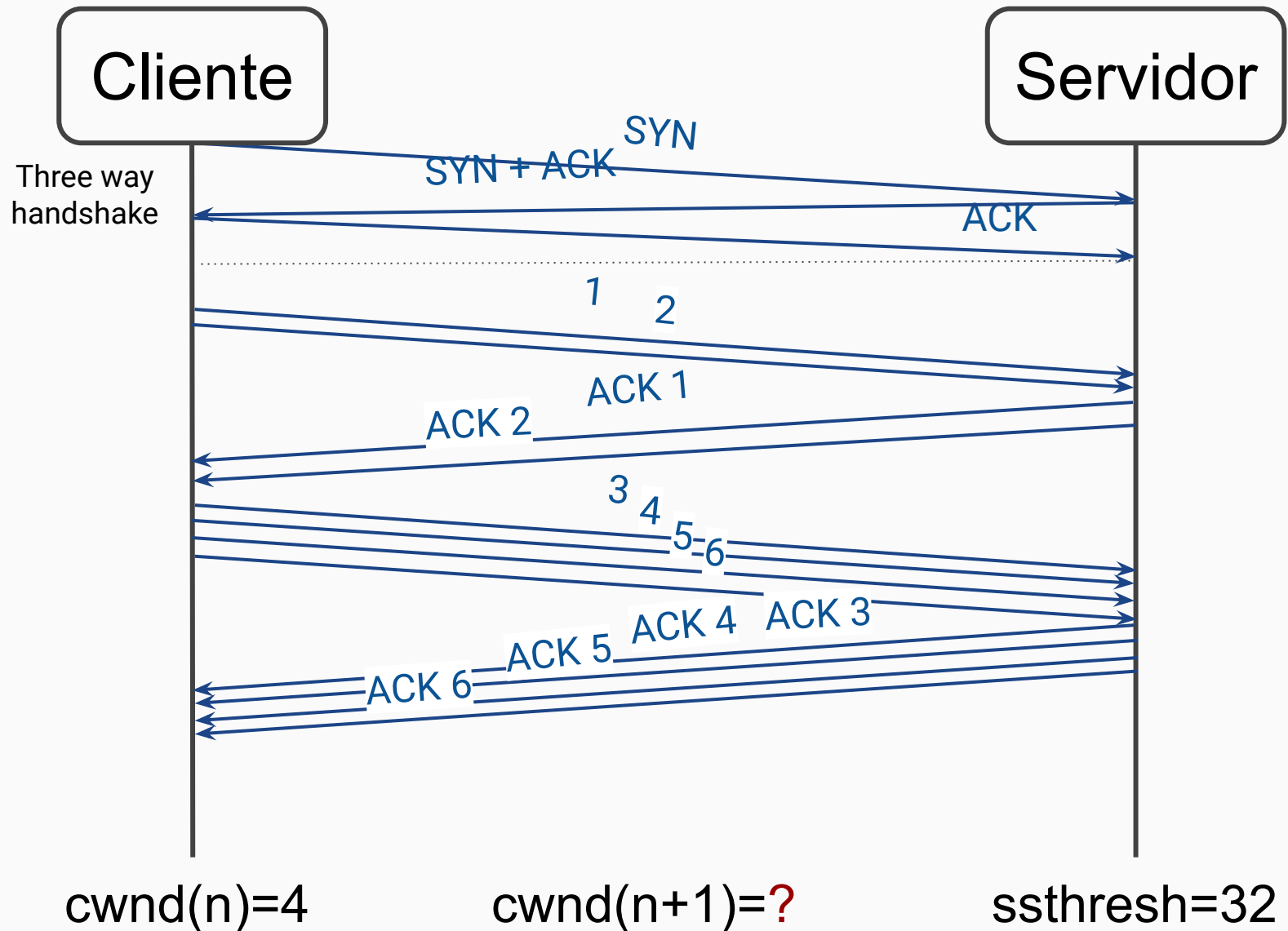
Resolución



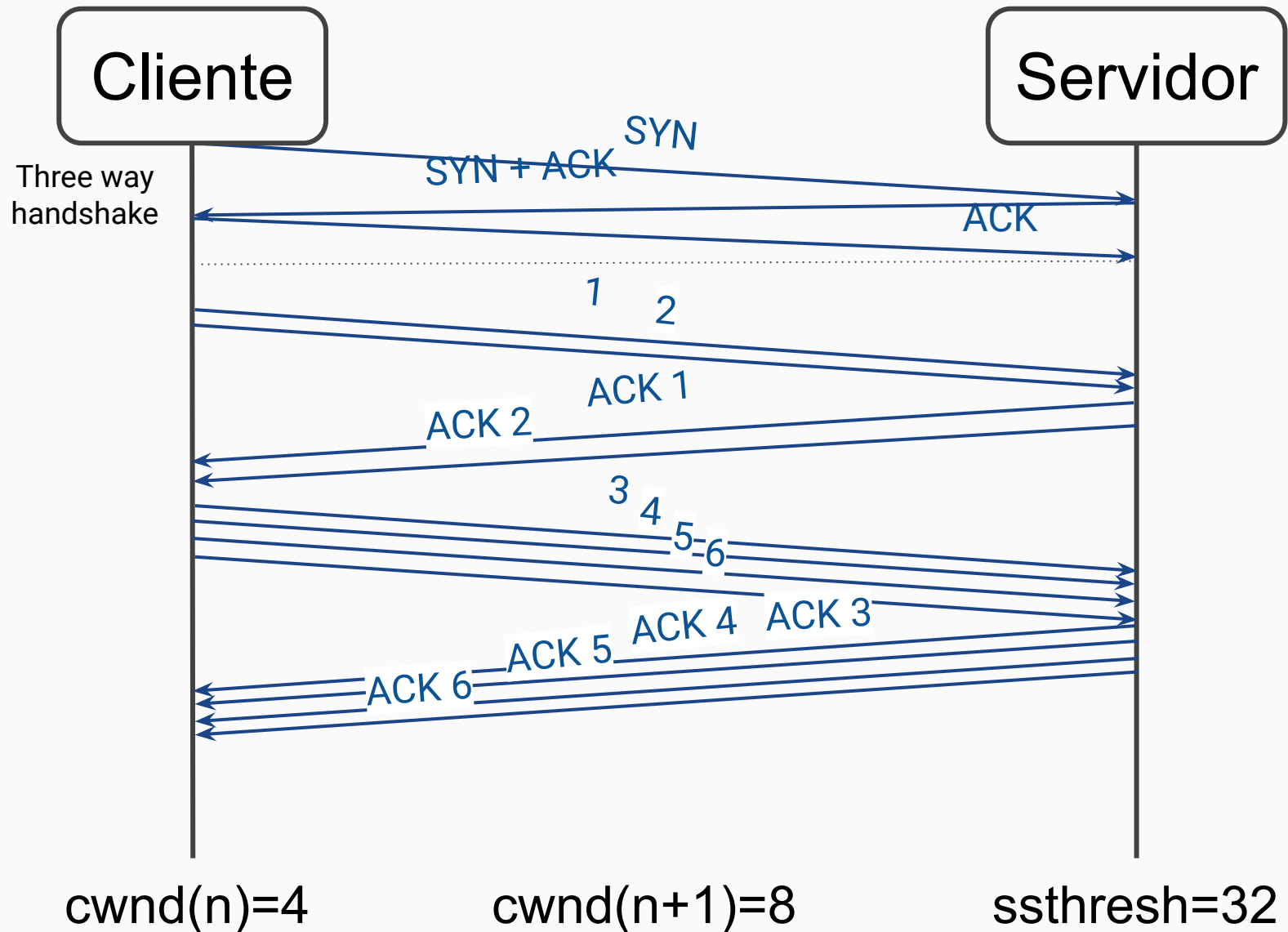
Resolución



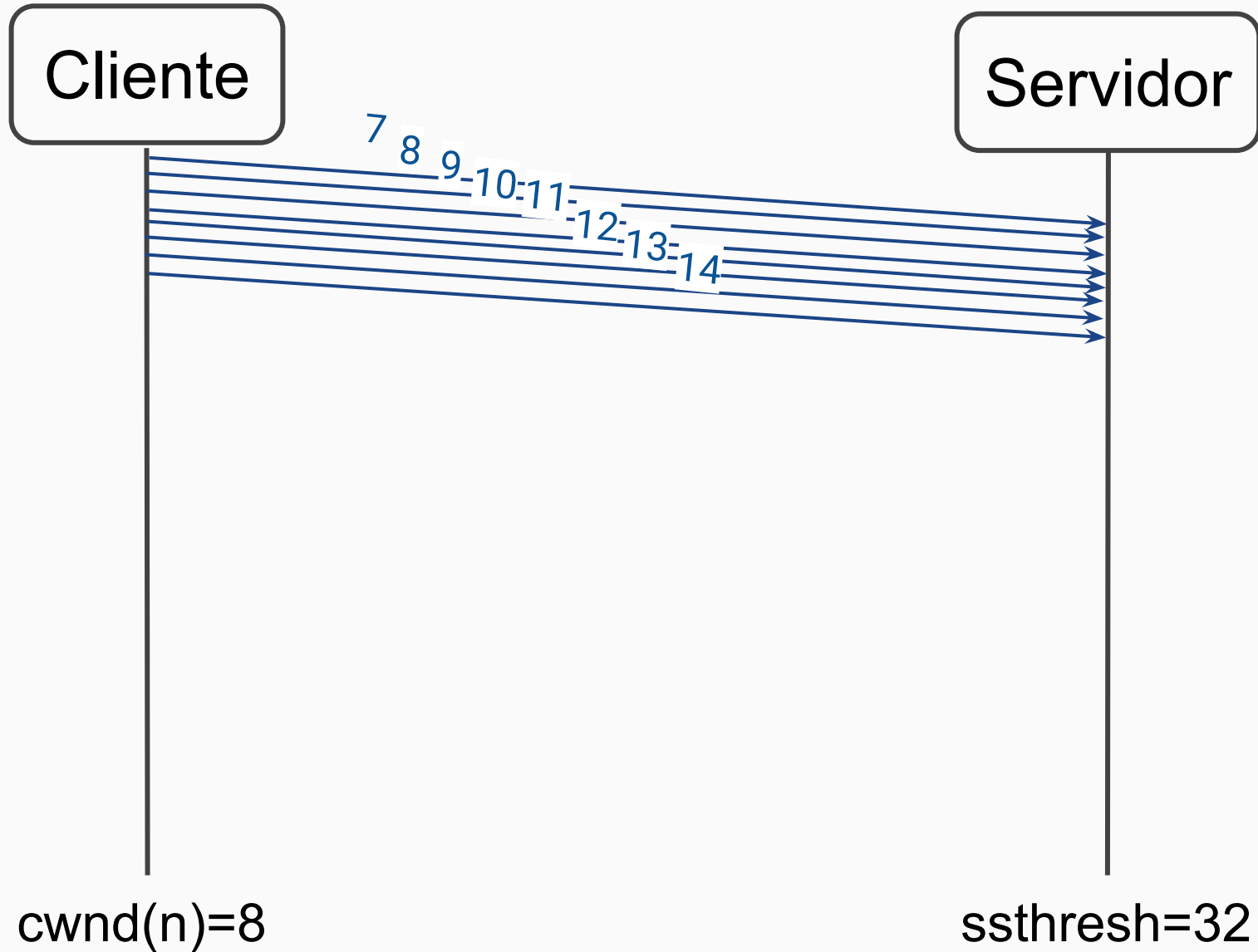
Resolución



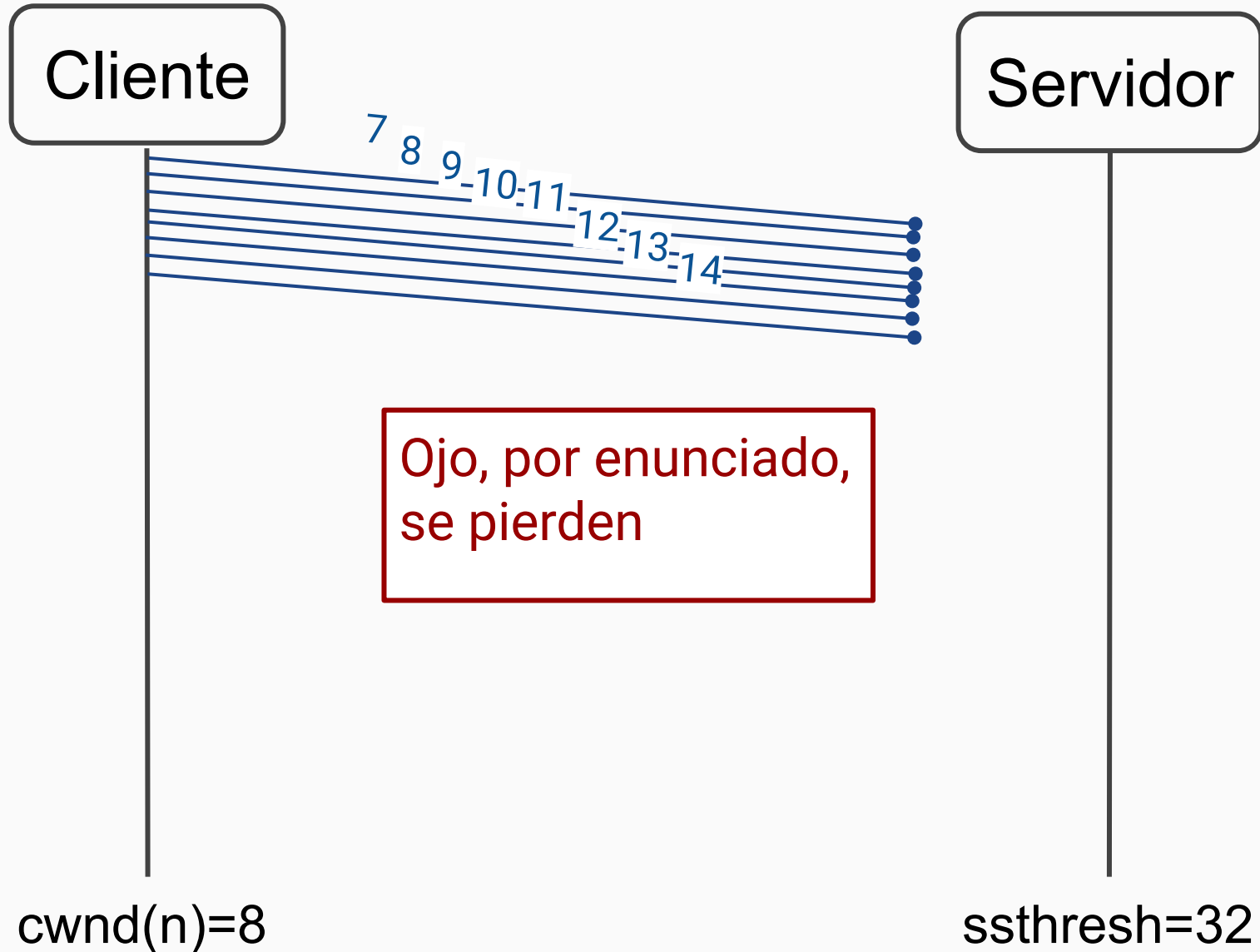
Resolución



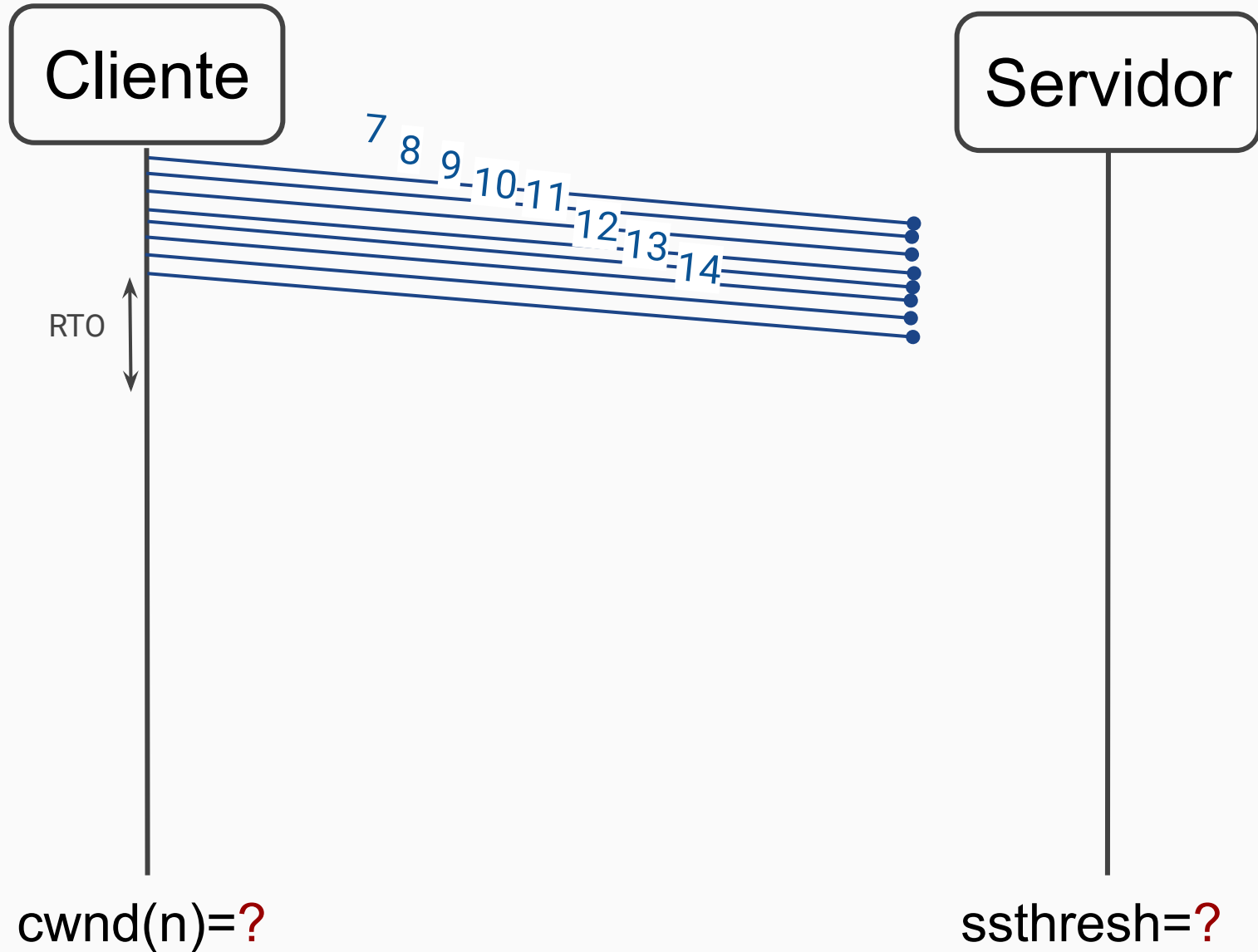
Slow start



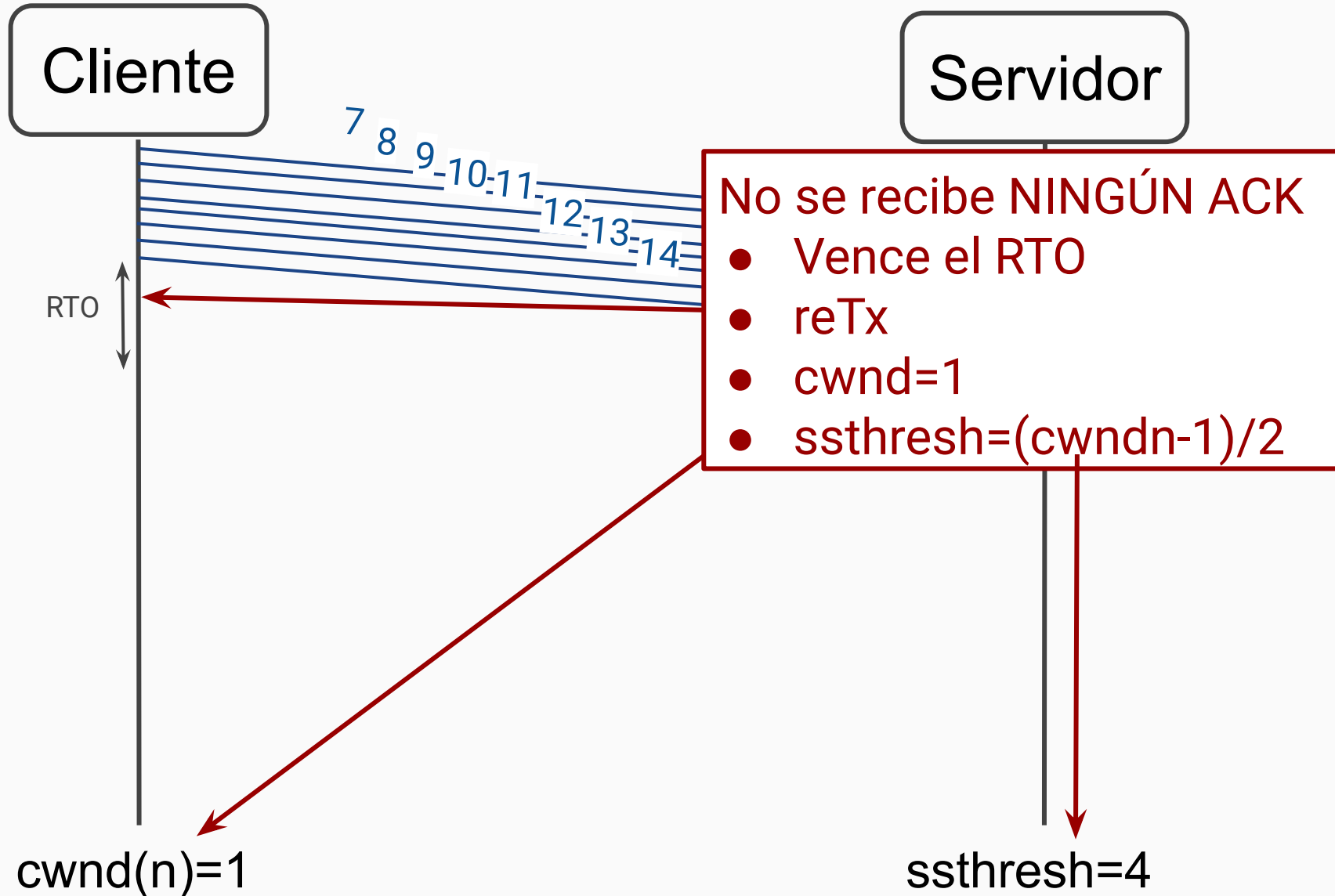
Slow start



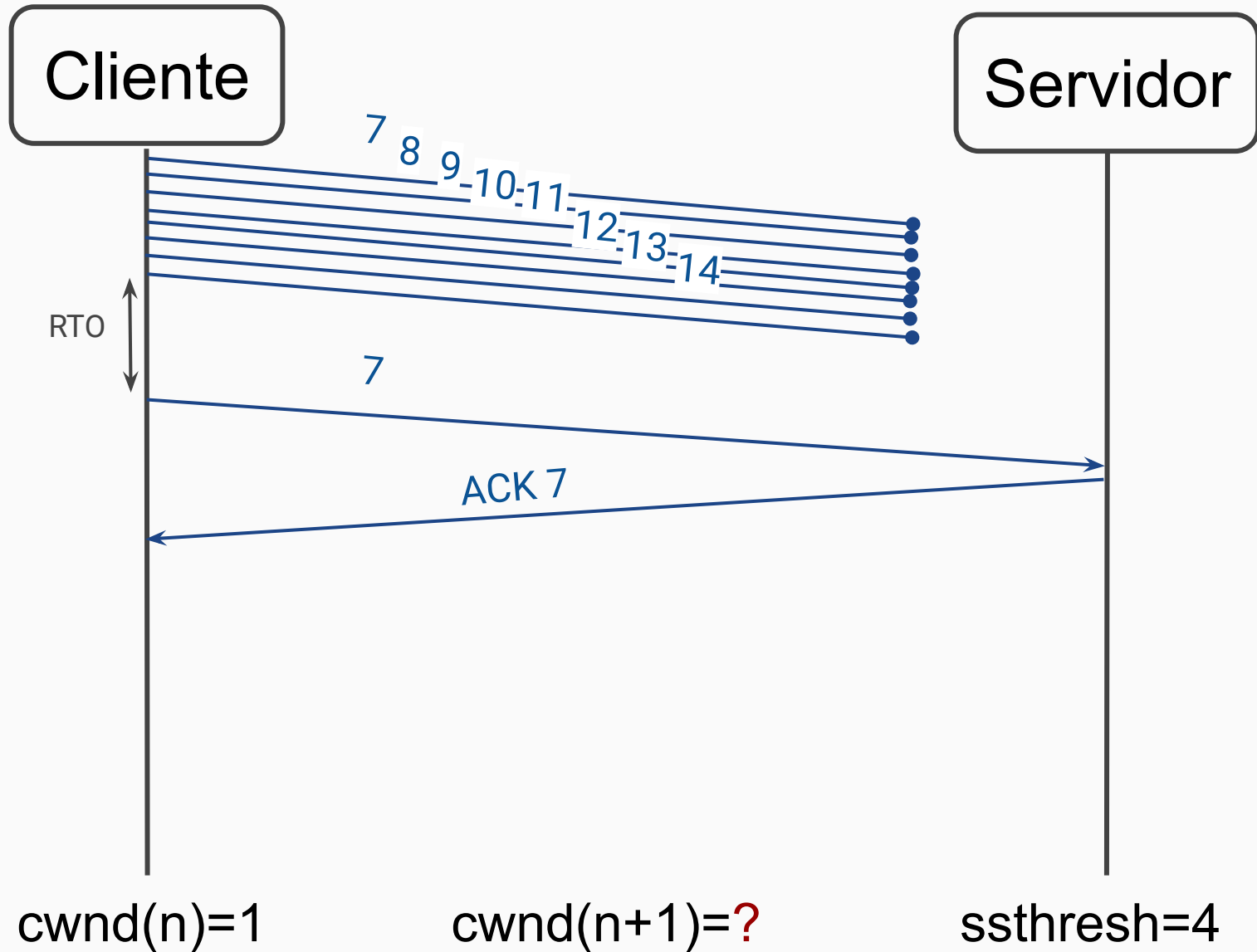
Slow start



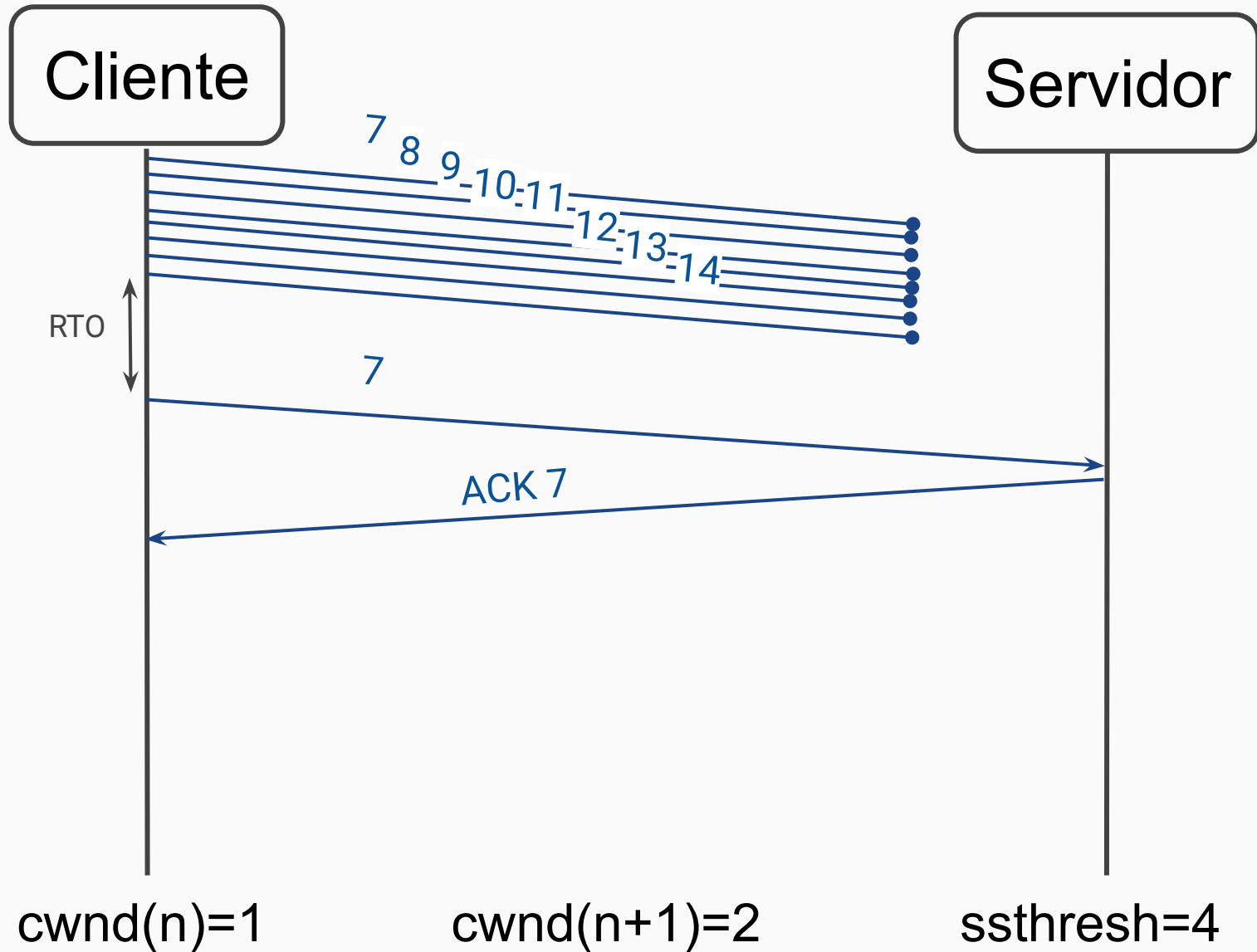
Slow start



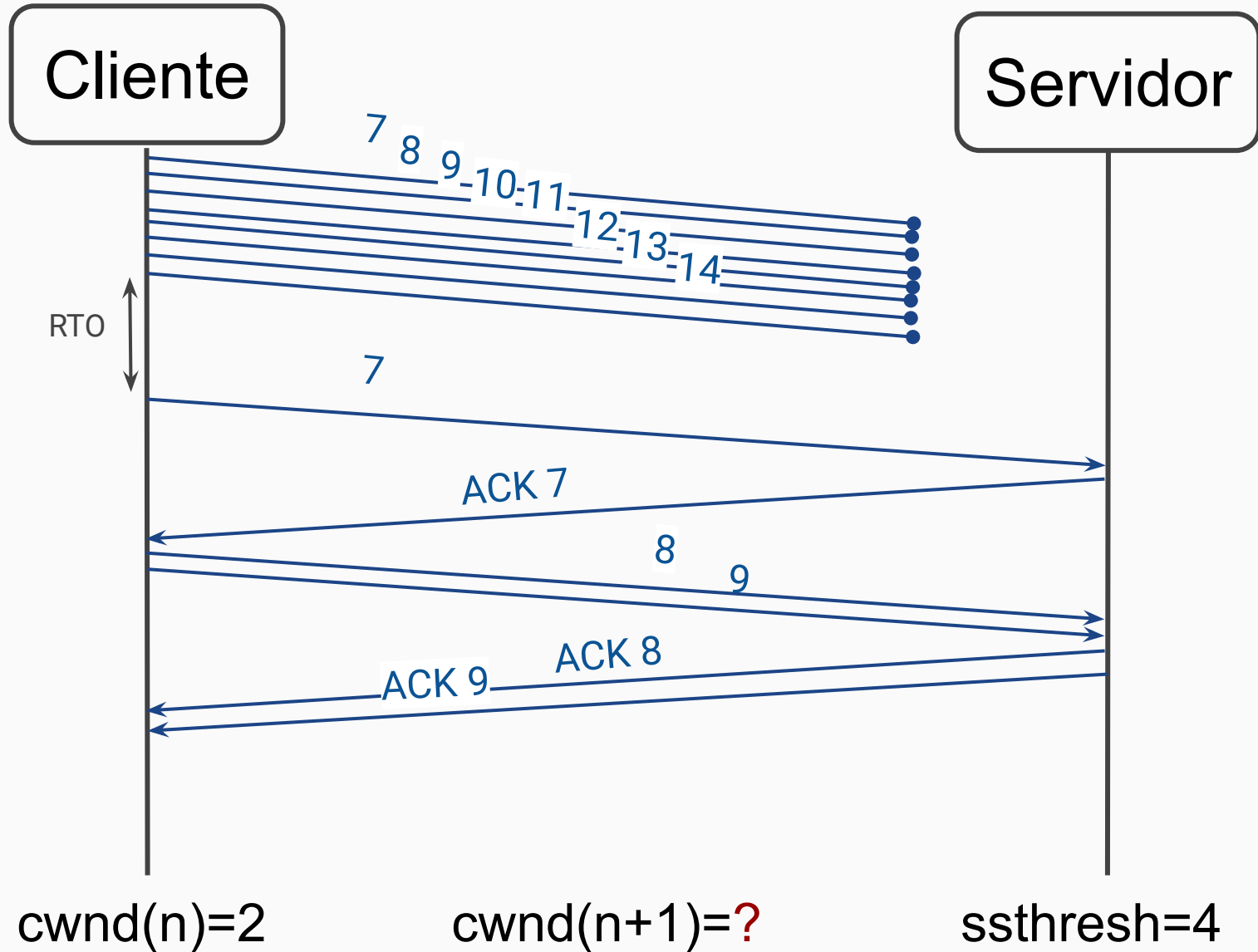
RTO Loss



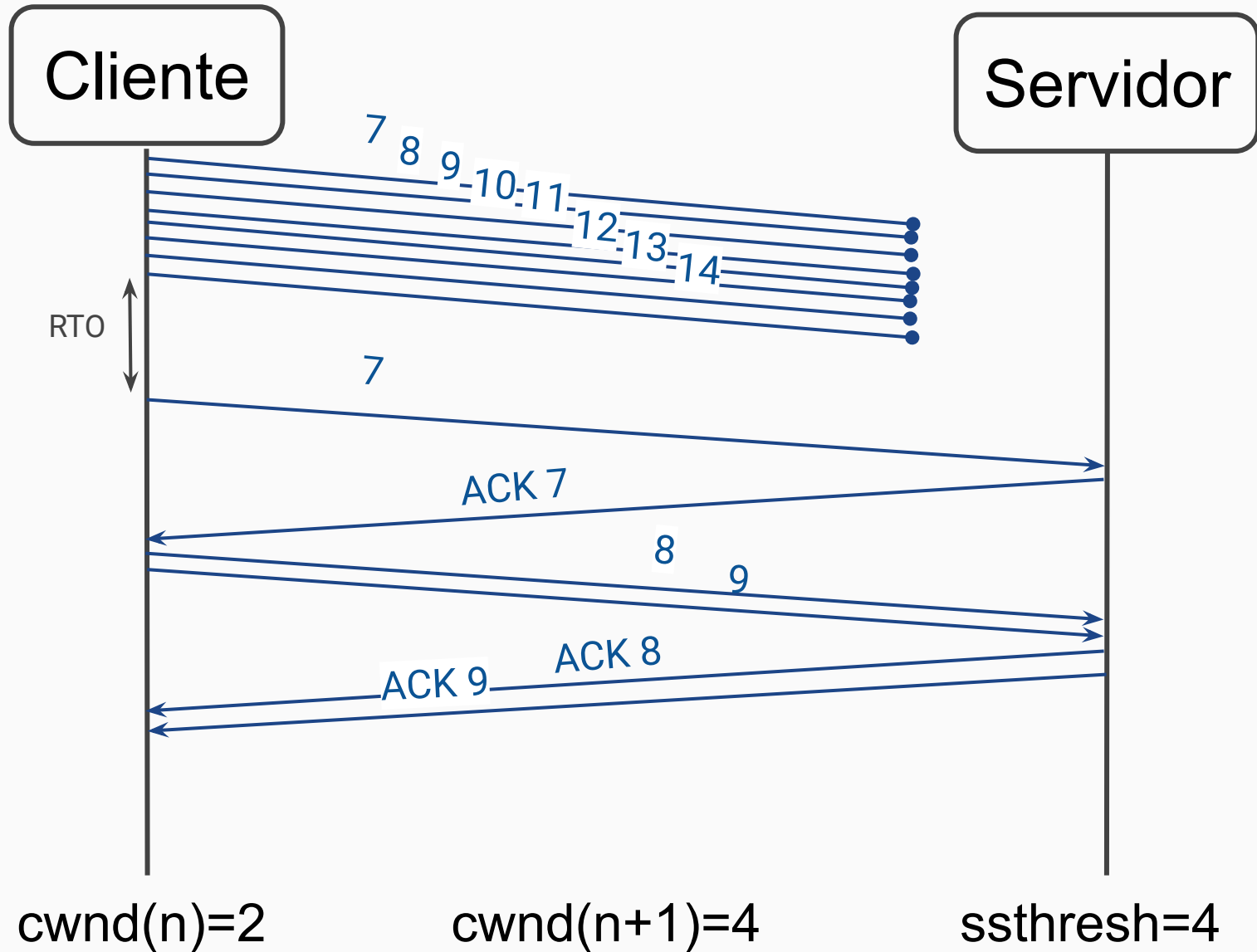
RTO Loss



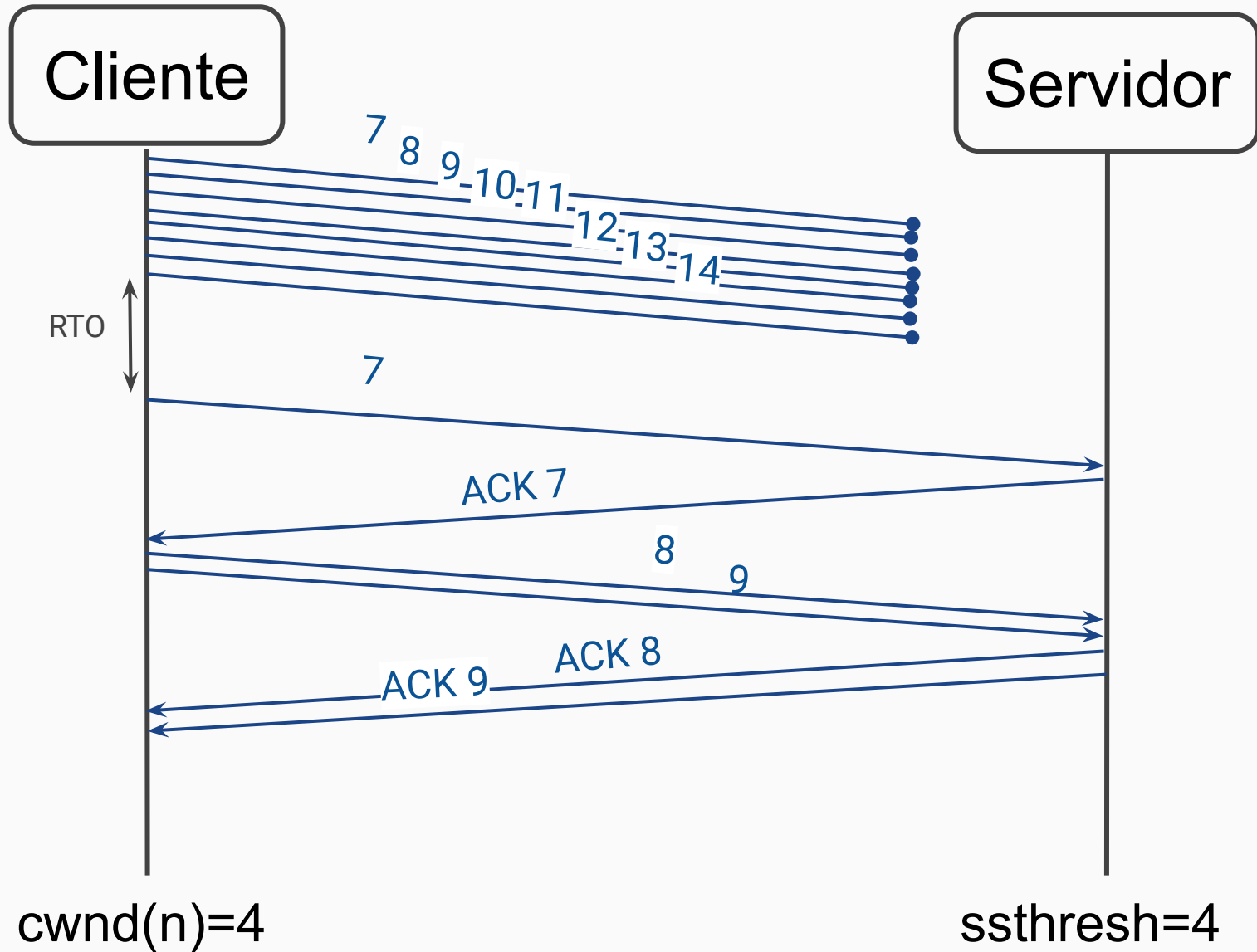
Slow start



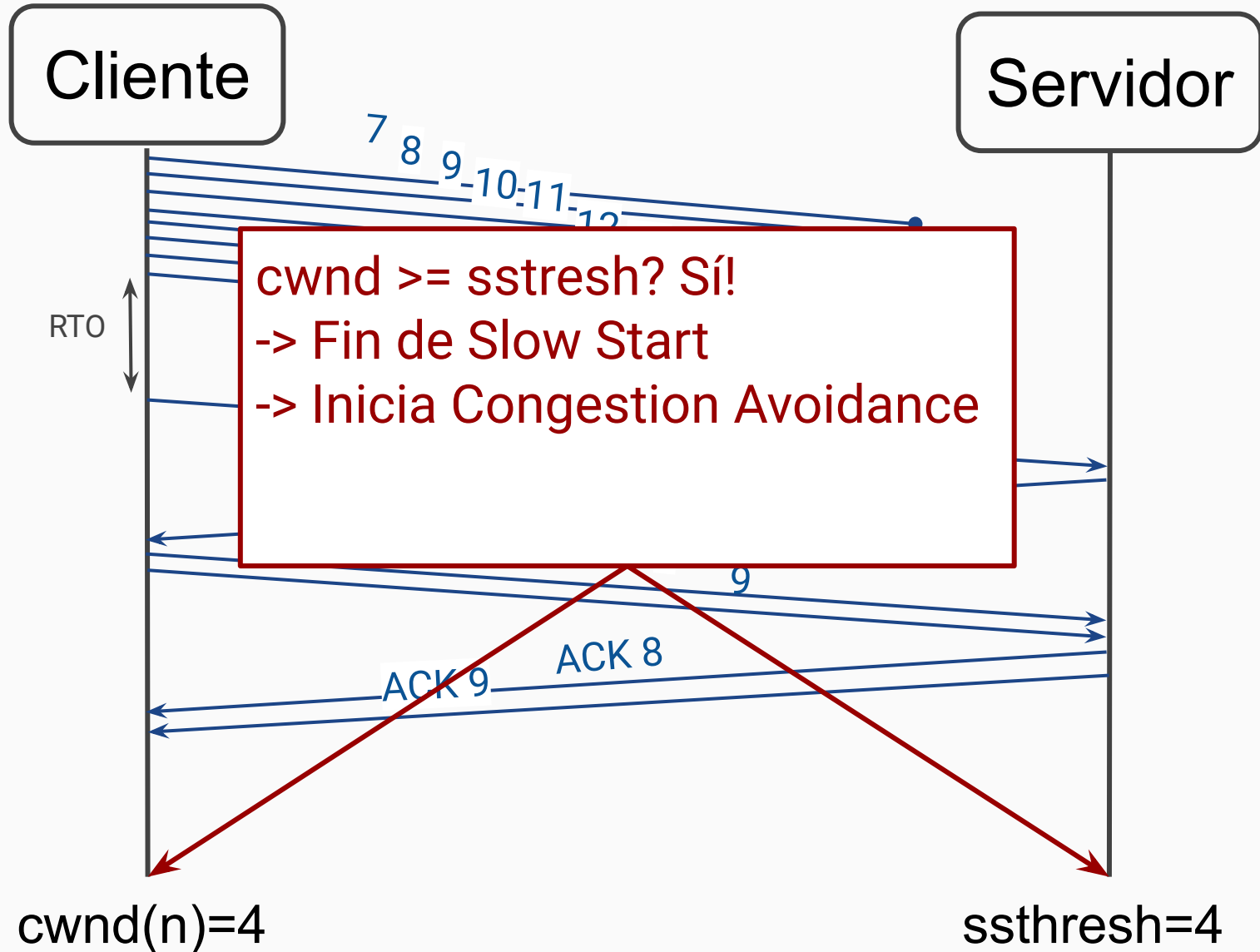
Slow start



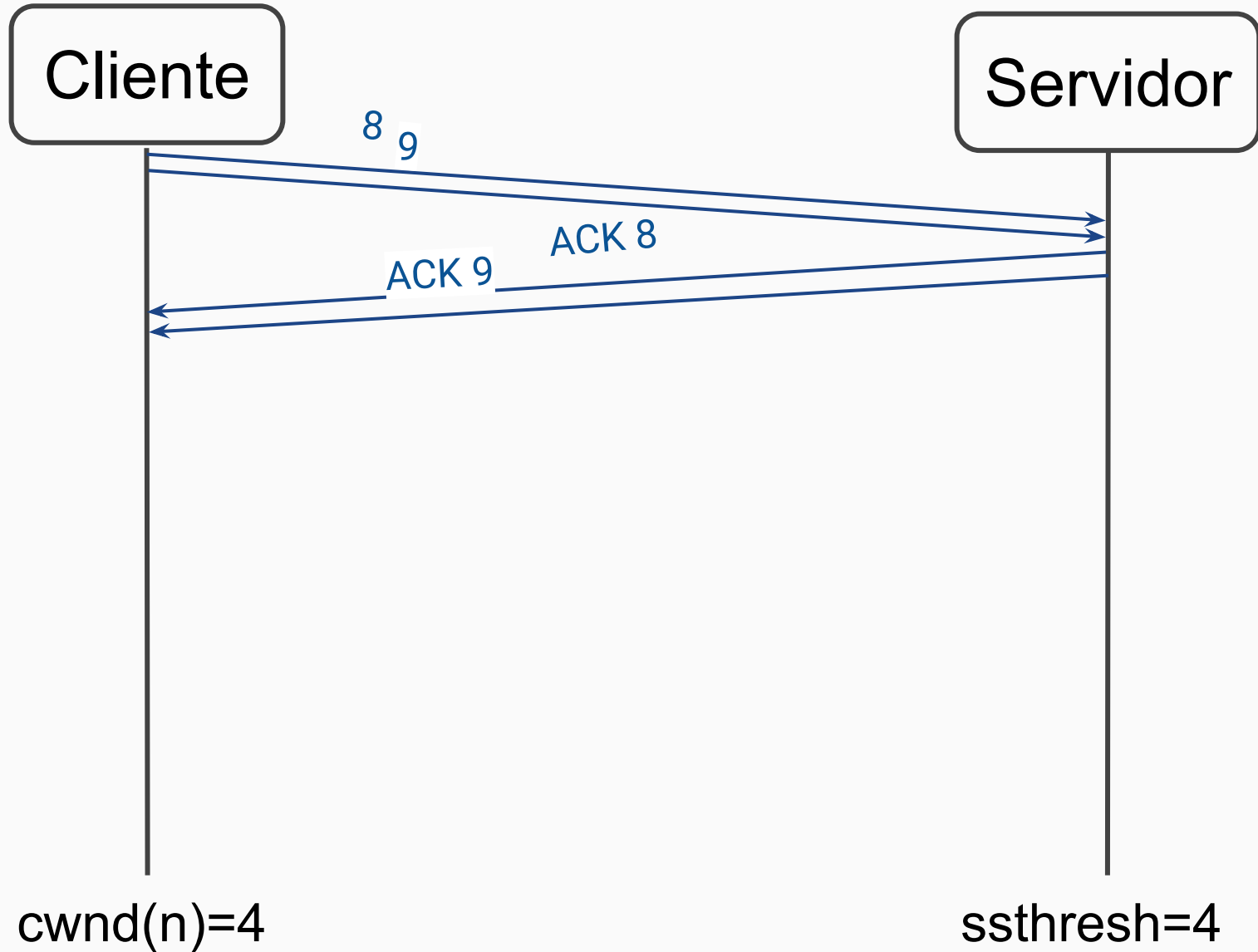
Slow start



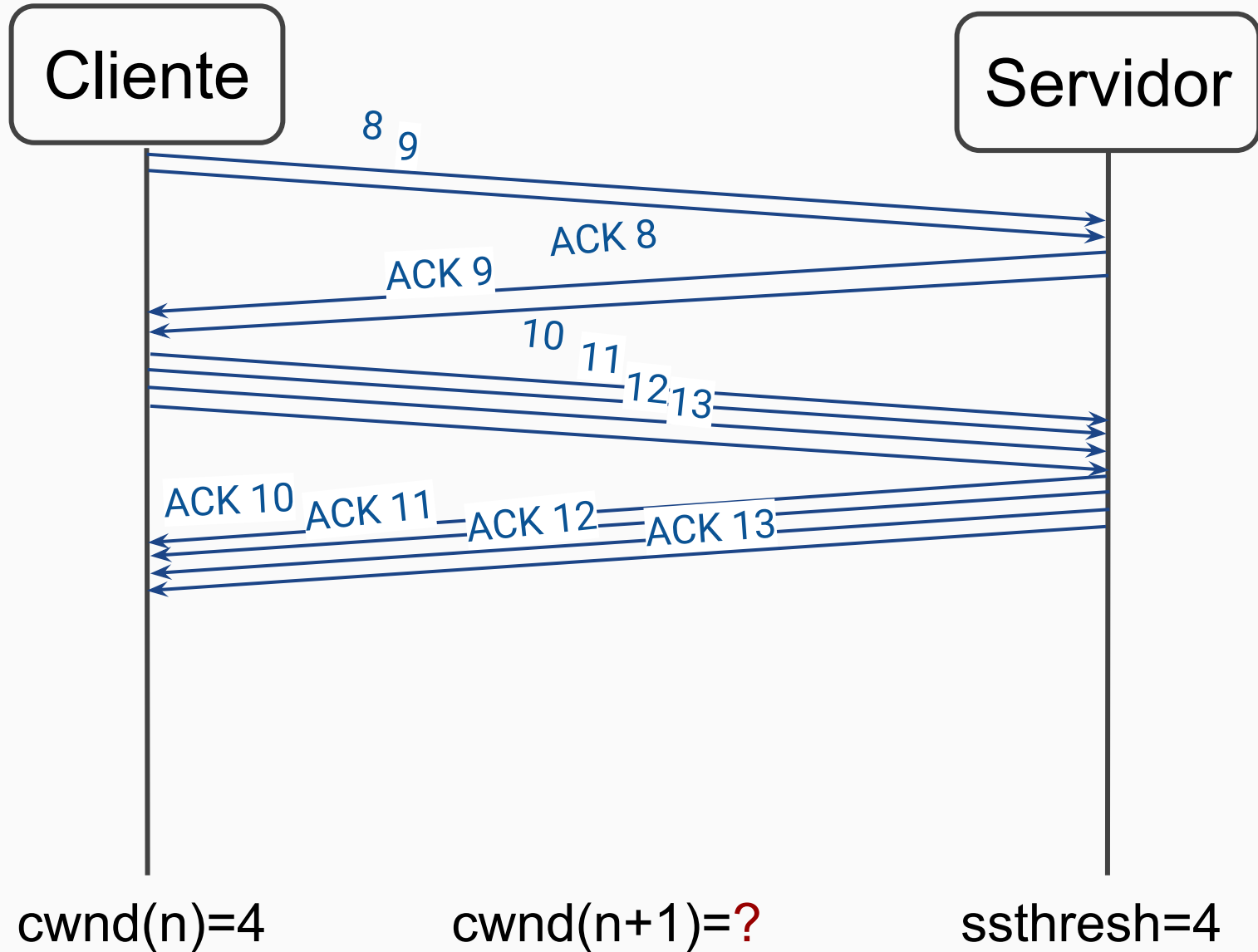
Slow start



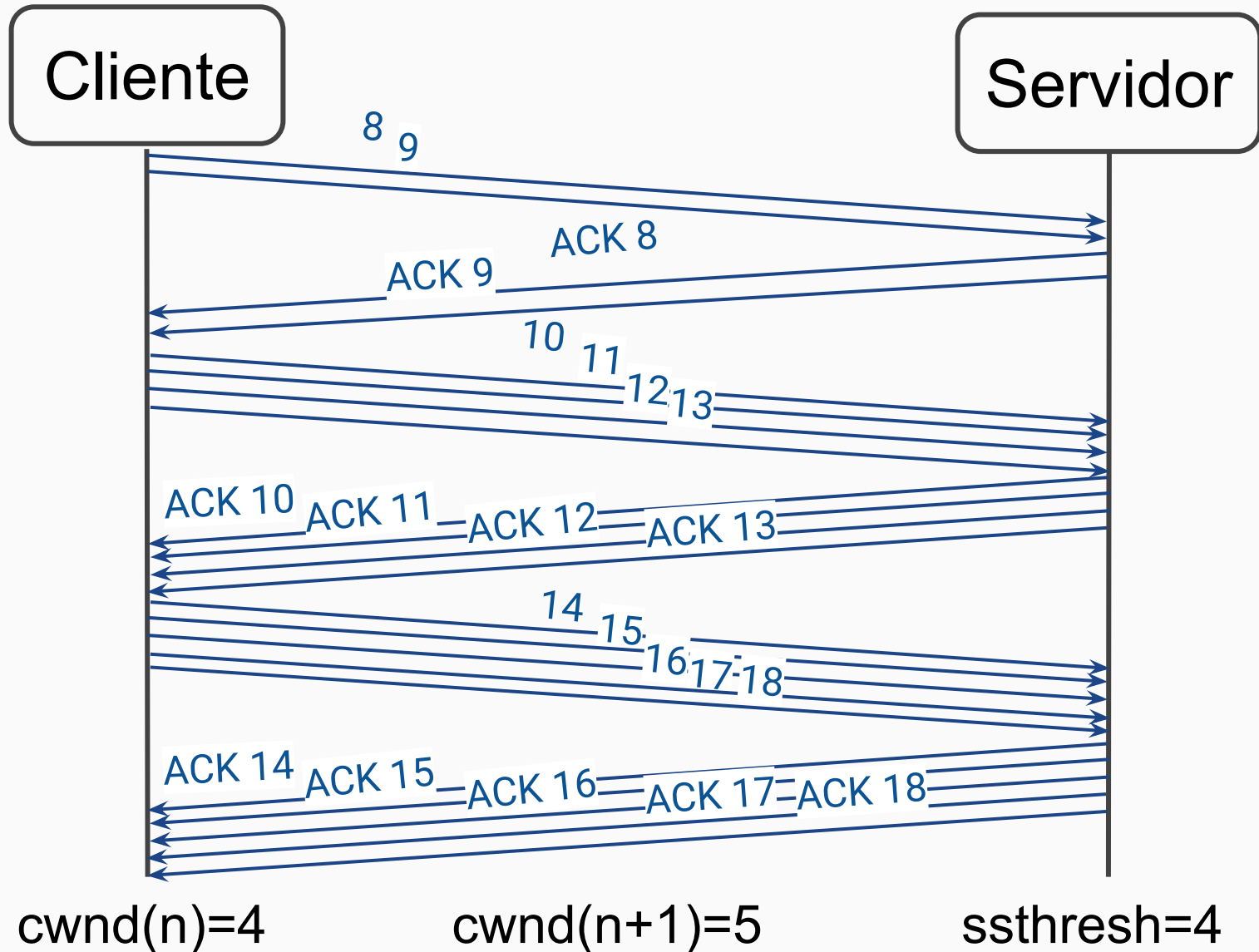
Slow start



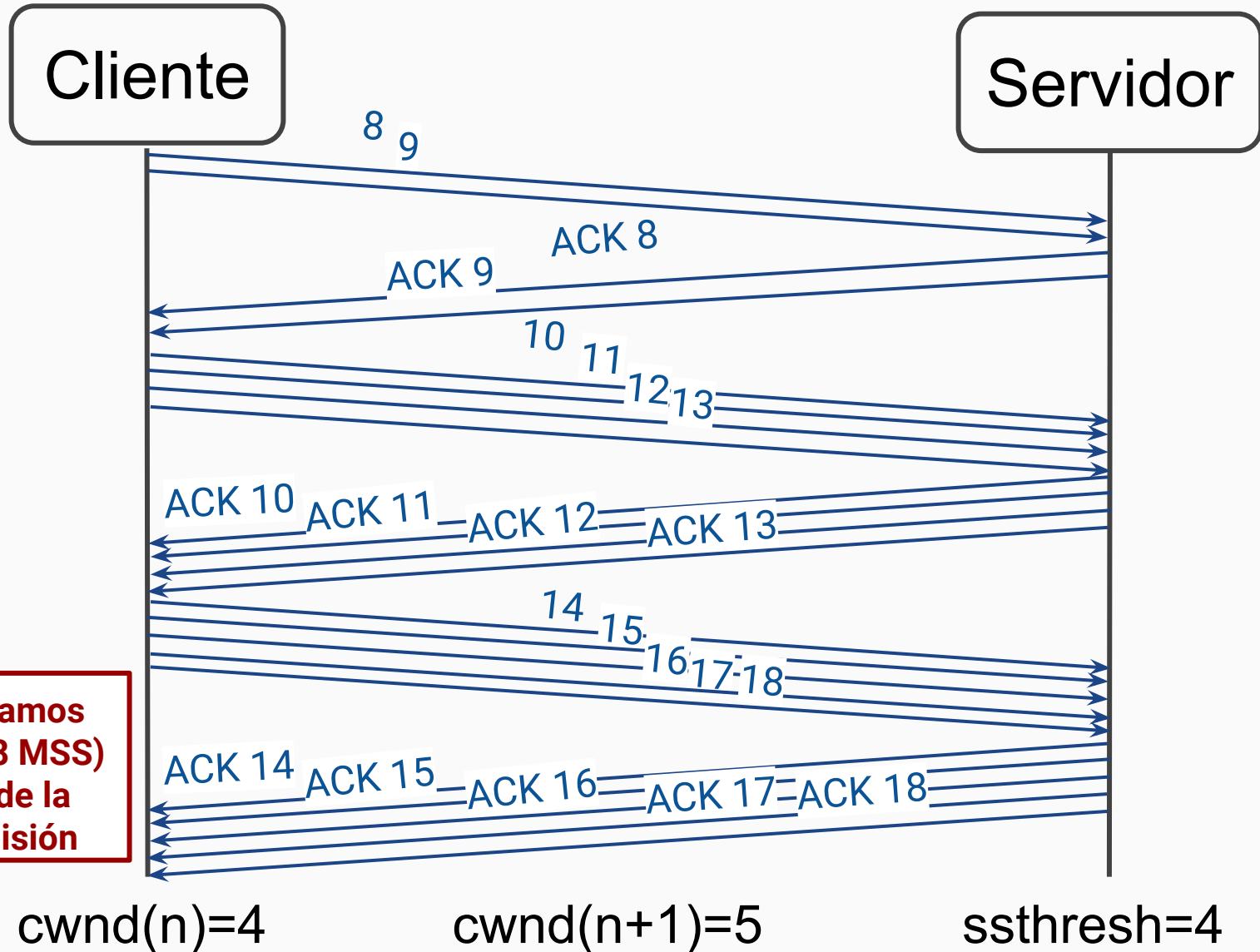
Slow start



Congestion avoidance



Congestion avoidance



Pérdida de paquetes

- ¿Qué pasaría si se pierde un paquete?
- Supongamos que se pierde **UN SOLO** segmento de la ráfaga en cuestión.
- Y además, que estamos usando el algoritmo de control de congestión de Tahoe

TCP : control de congestión

ACK duplicados (llegan algunos paquetes)

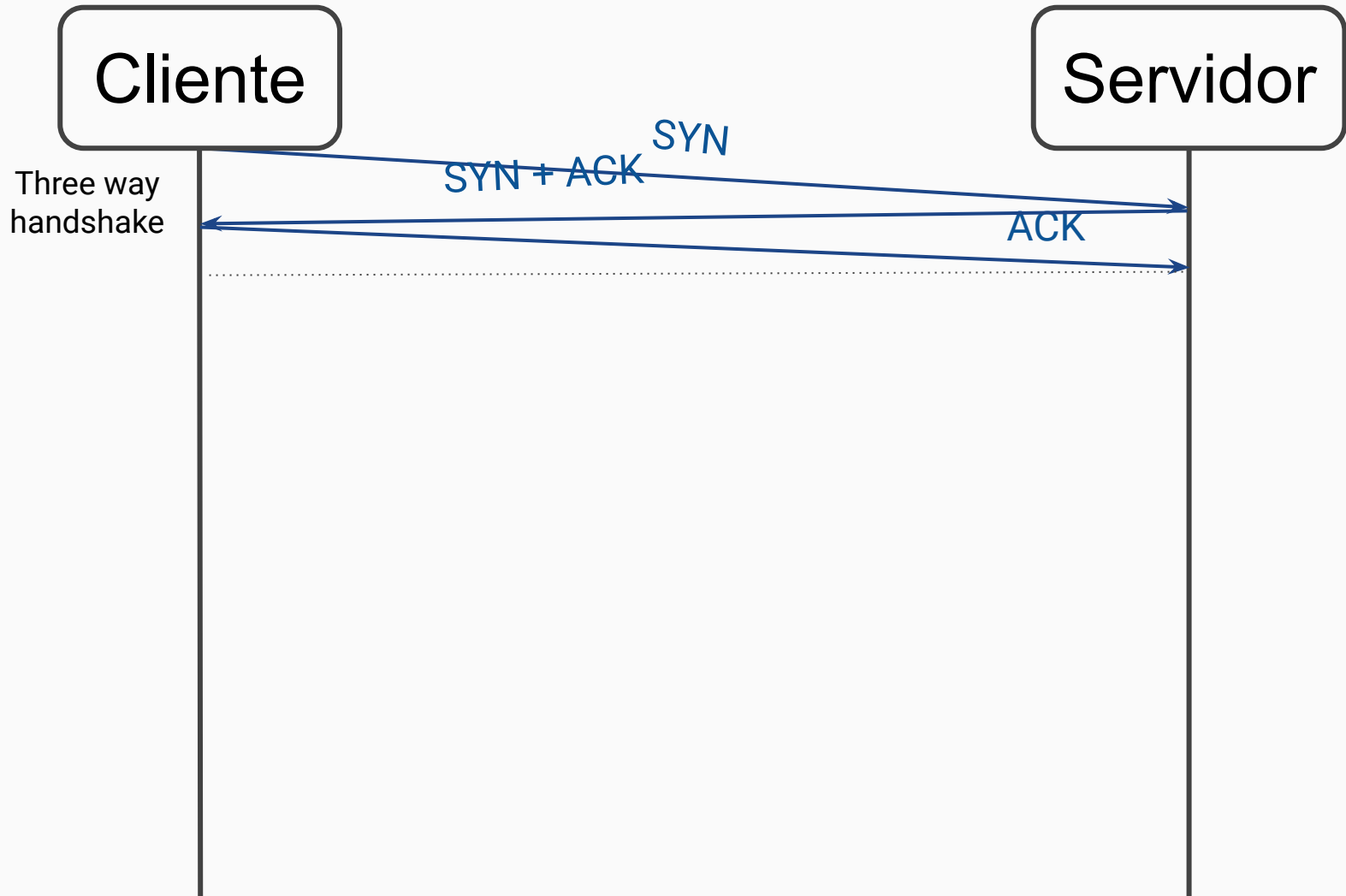
- Ocurre cuando se reciben 4 ACKs para el mismo segmento (3 dupACKs)
- Se inicia **Fast Retransmit**
- Depende del algoritmo implementado
 - **Fast Recovery**
 - **Slow Start**

TCP : control de congestión

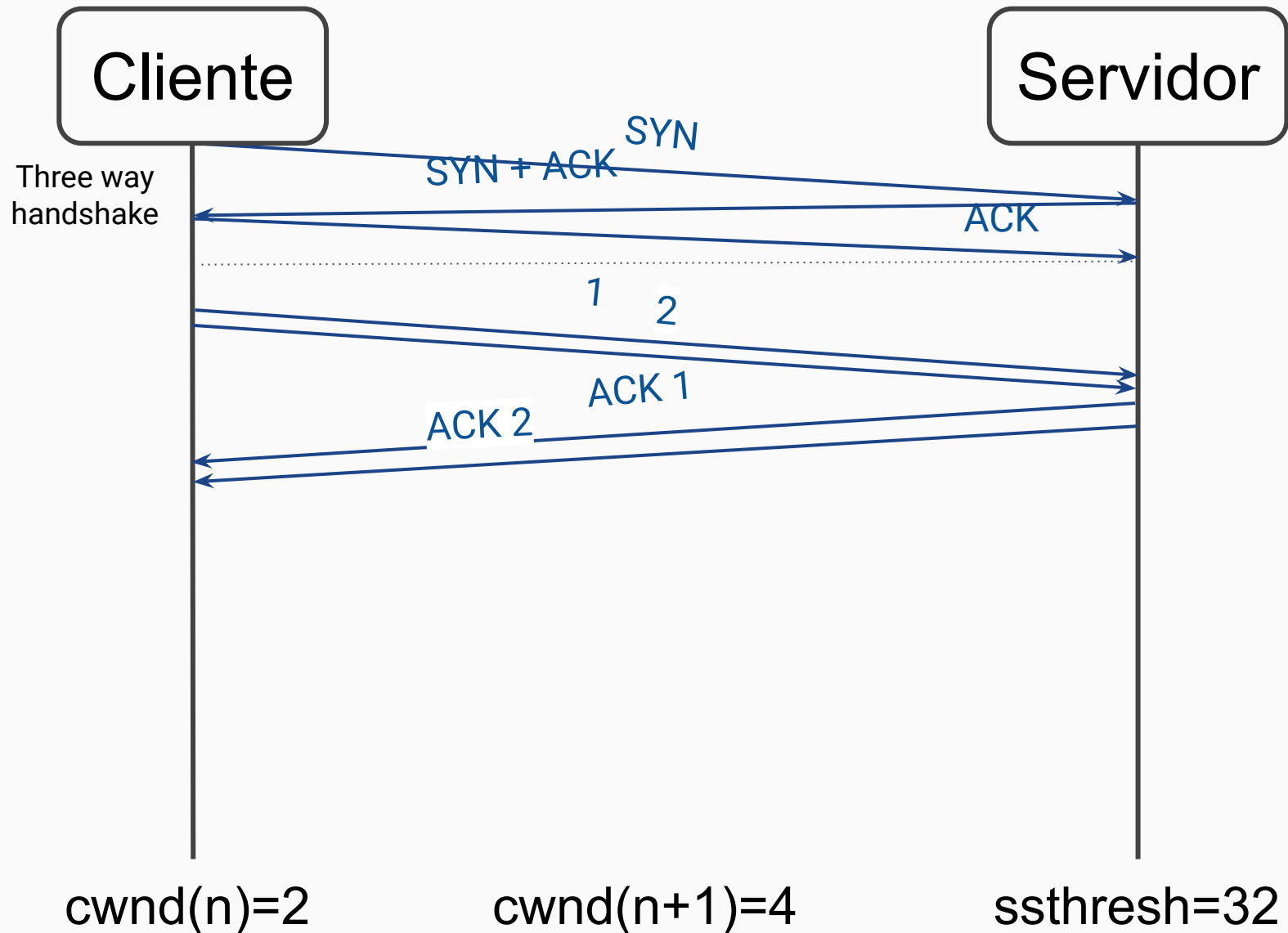
Algoritmo TCP Tahoe

- Le da al caso de los ACKs duplicados el mismo tratamiento que a un RTO (**Fast Retransmit** y luego **Slow Start**)
- Ver ejemplo

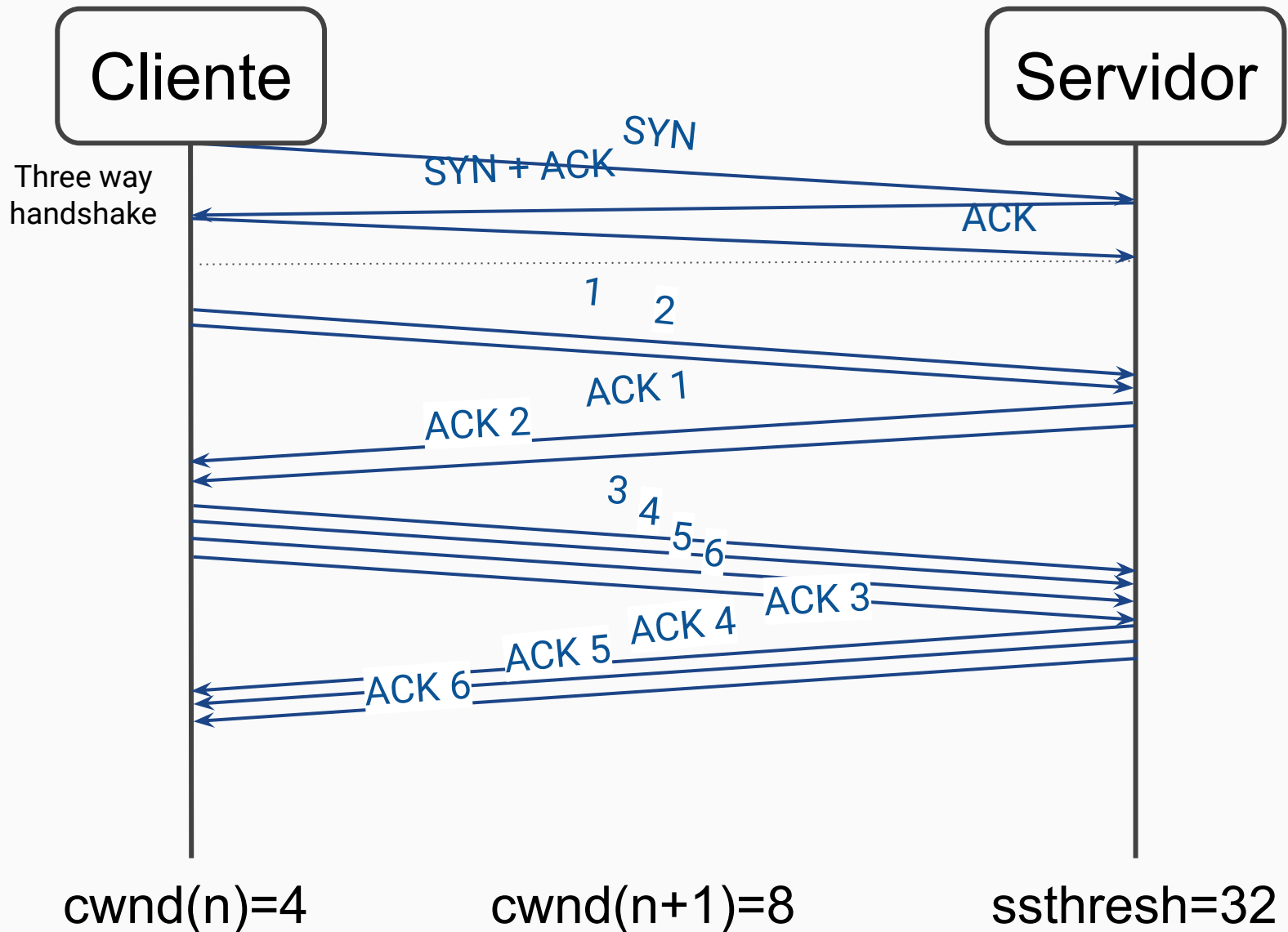
Three way handshake



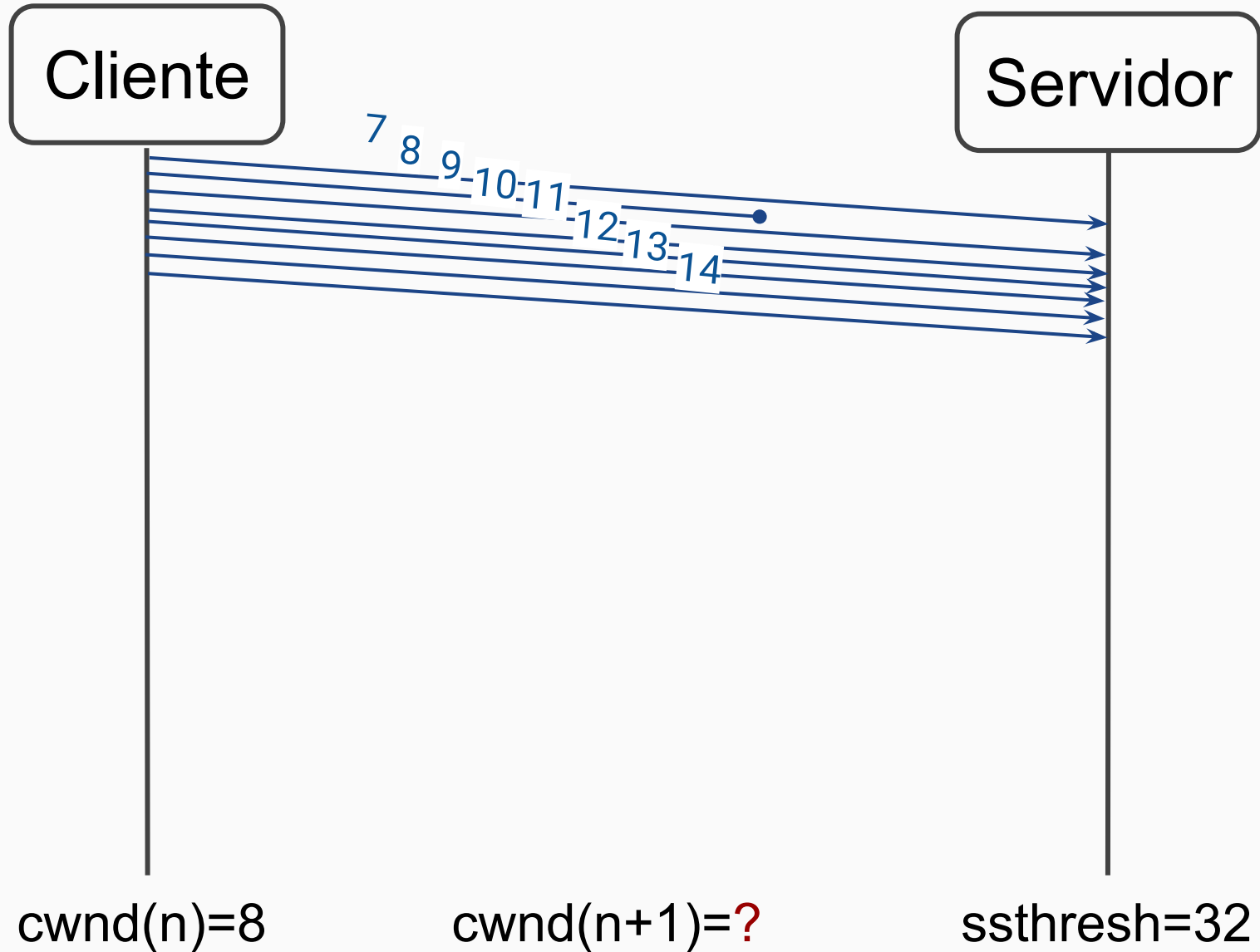
Slow start



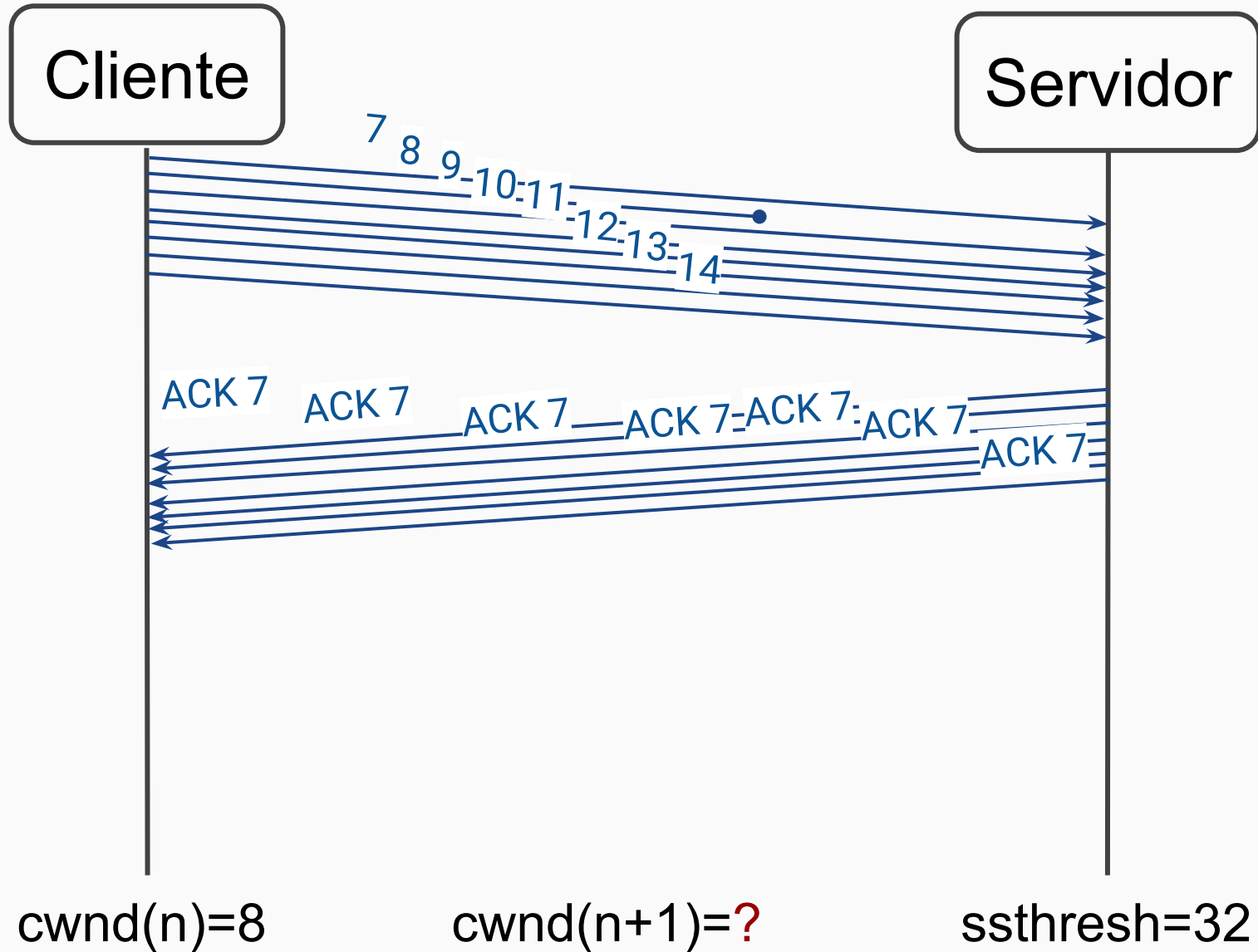
Slow start



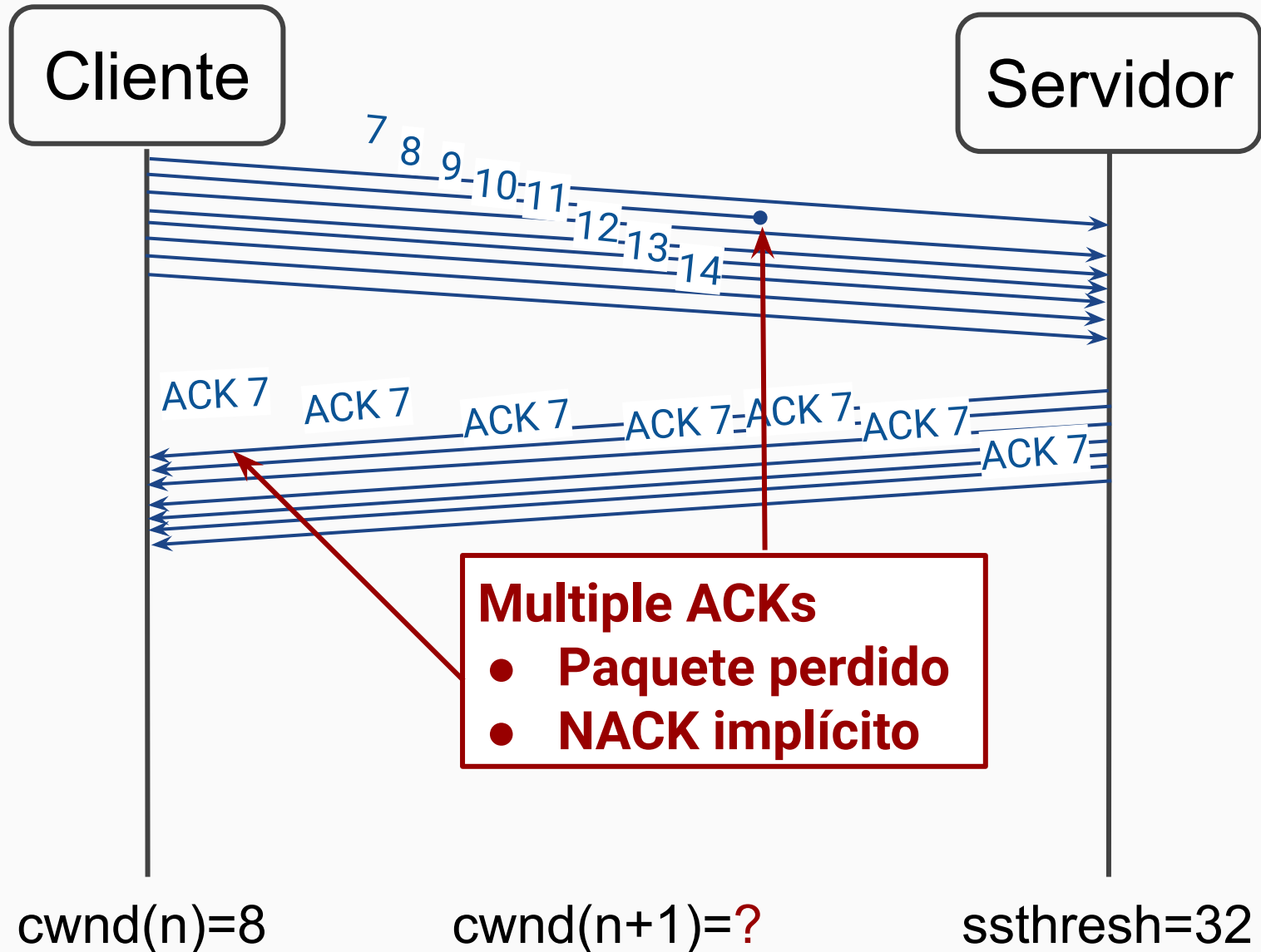
Slow start



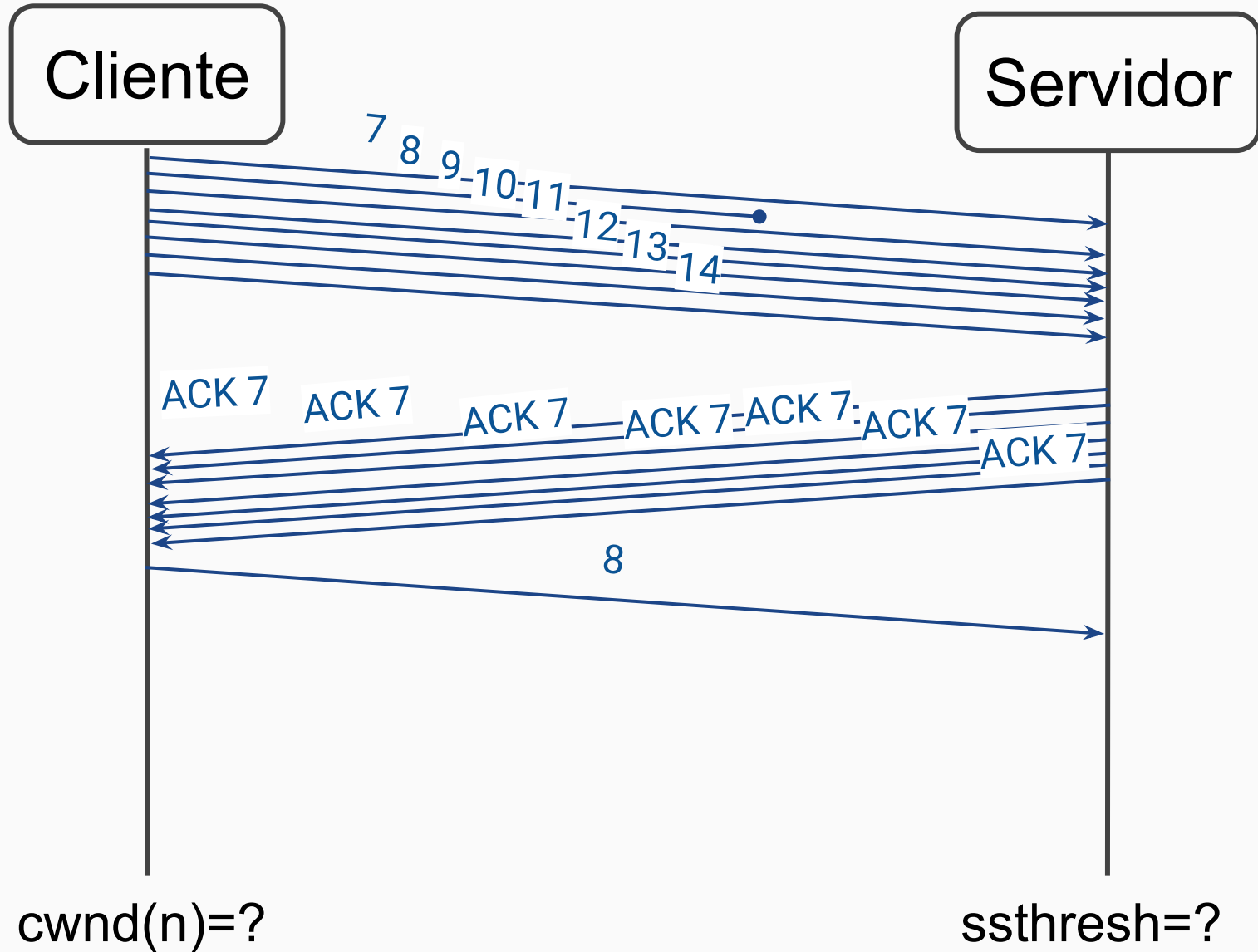
Slow start



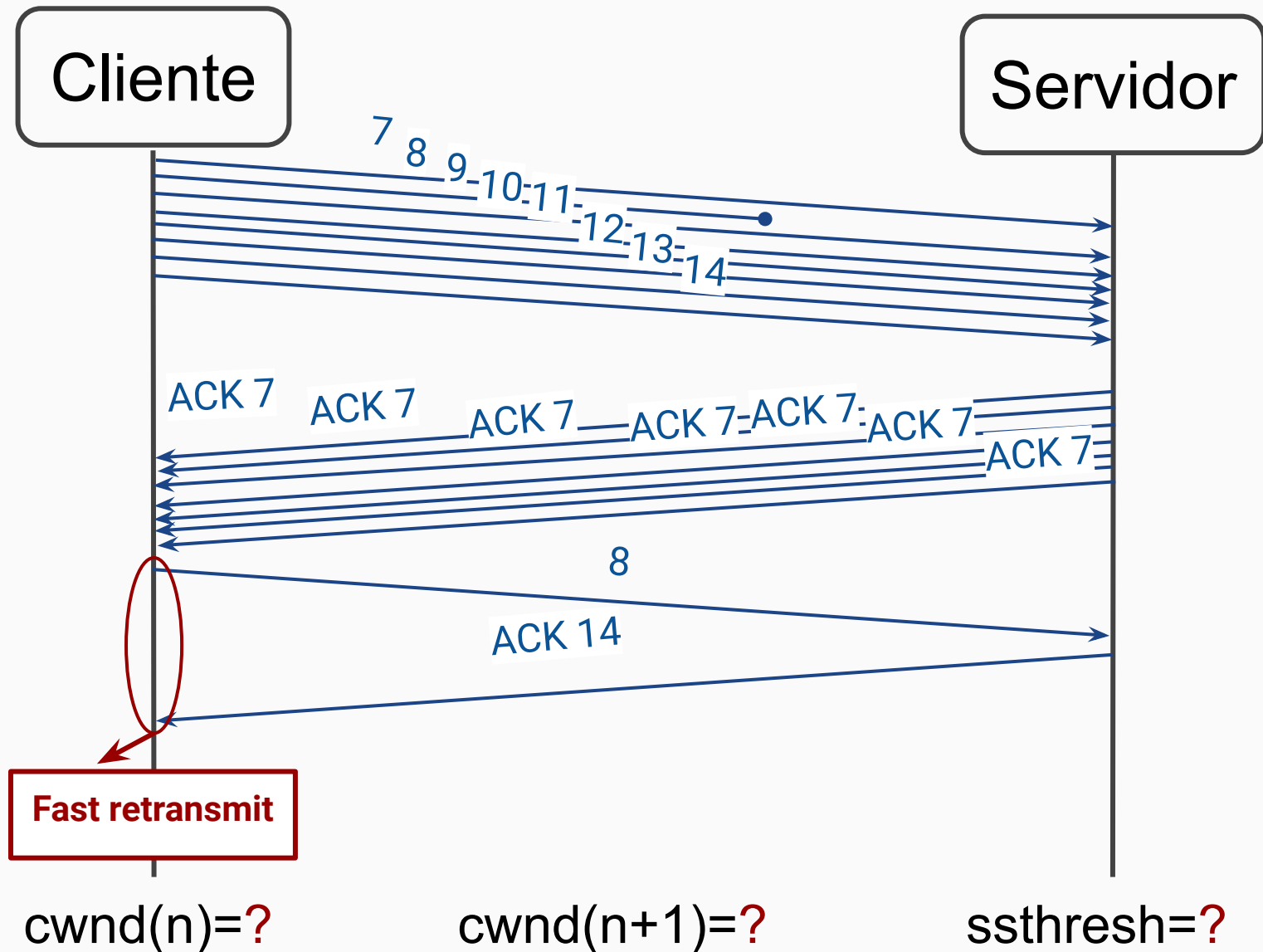
Slow start



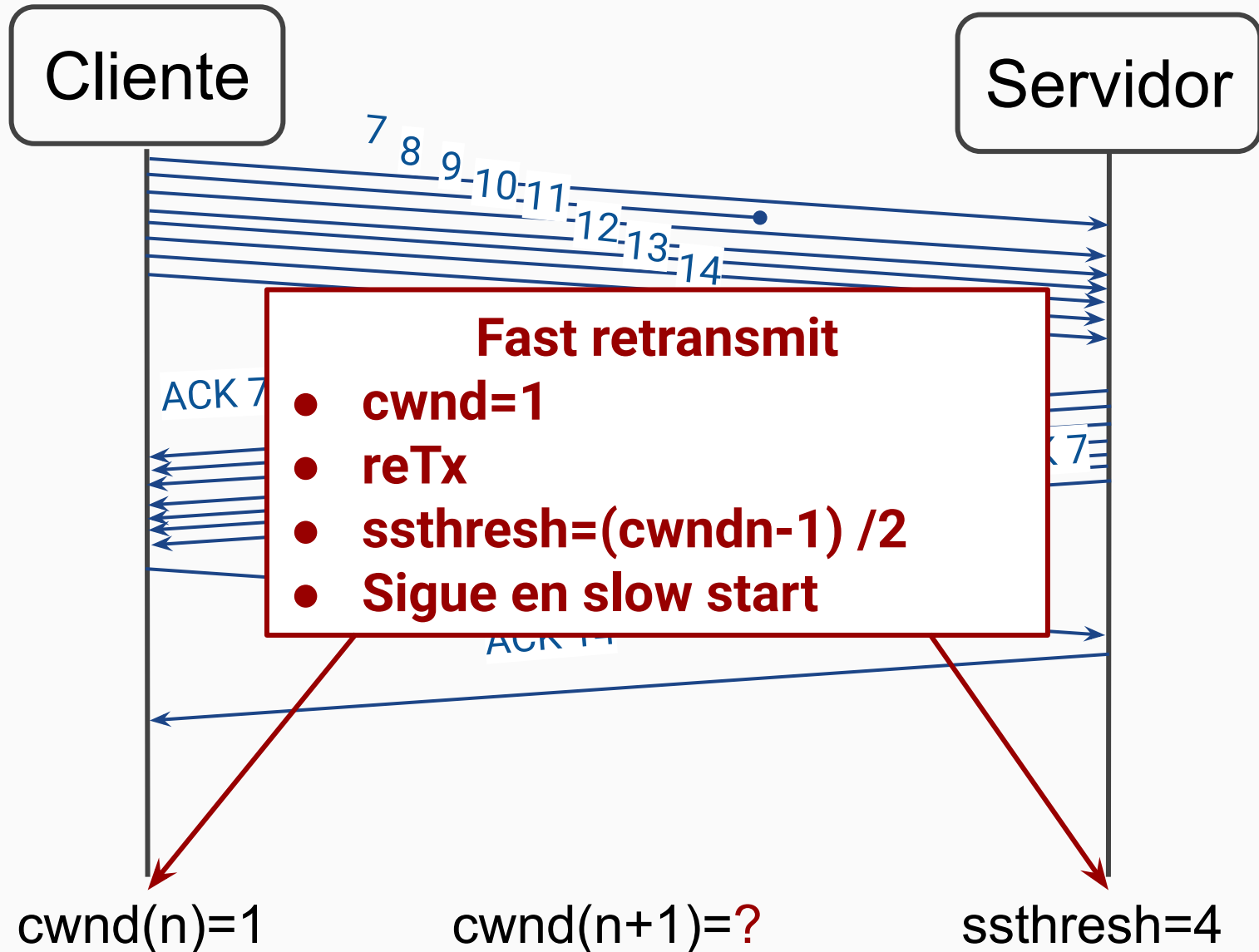
Slow start



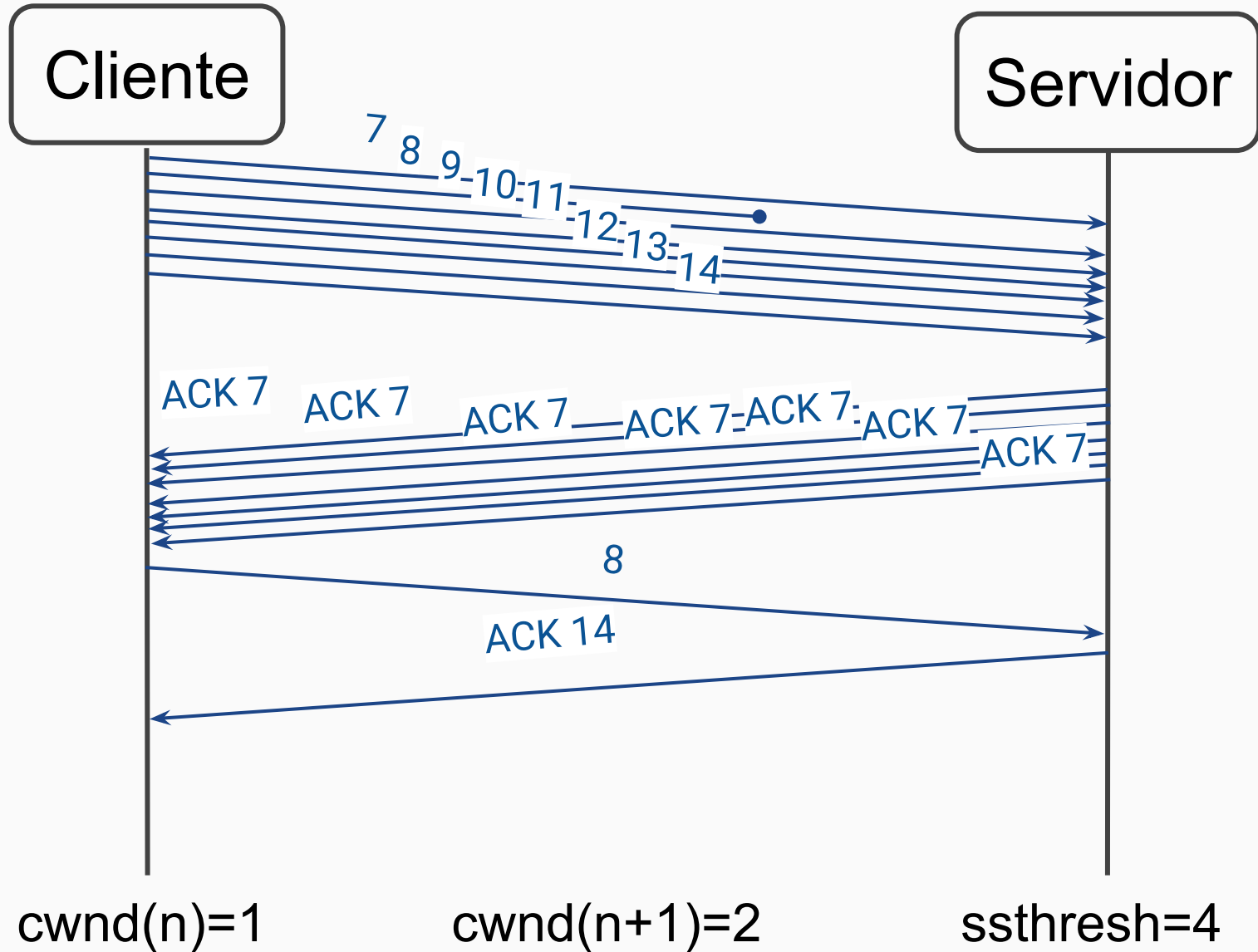
Fast Retransmit



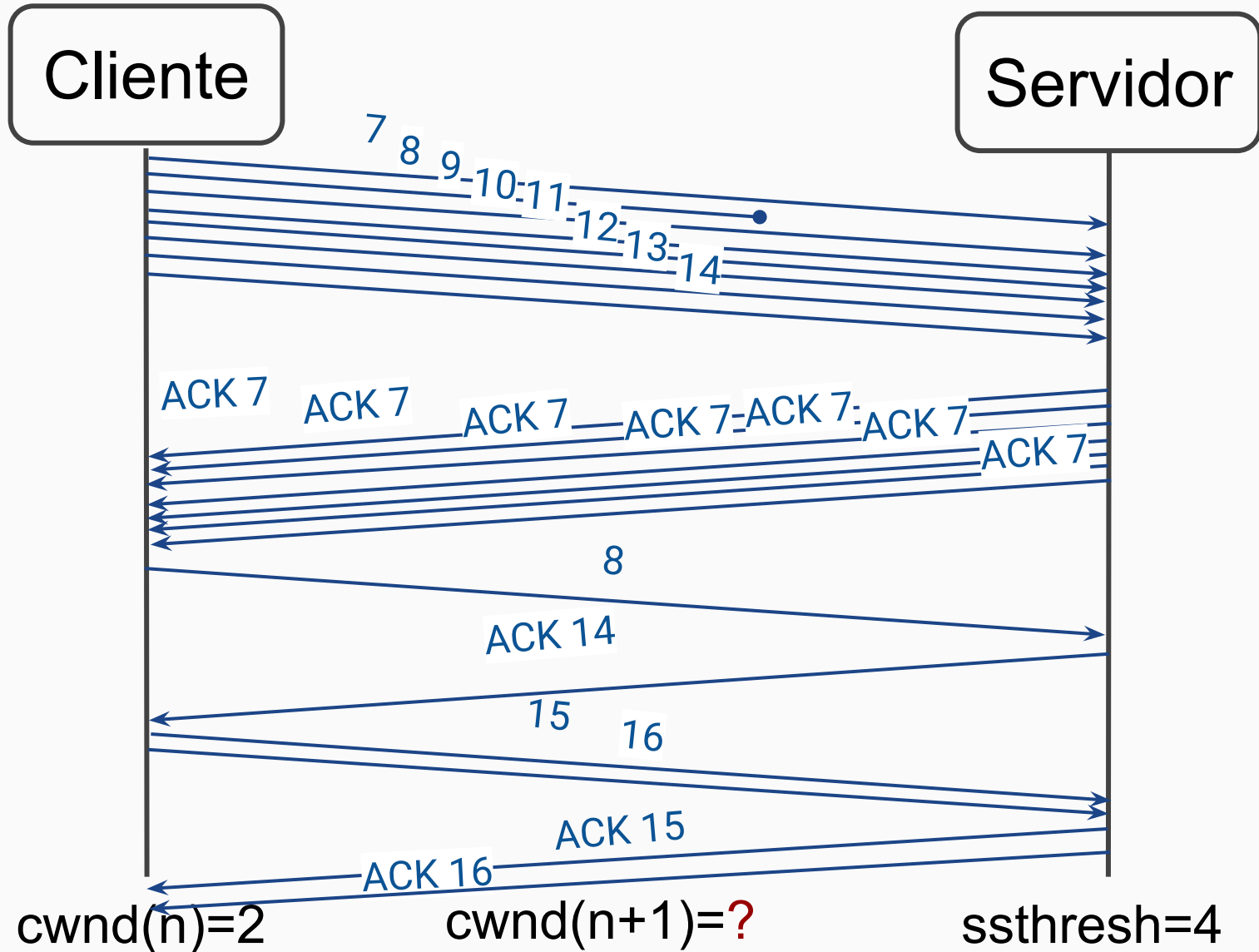
Fast Retransmit



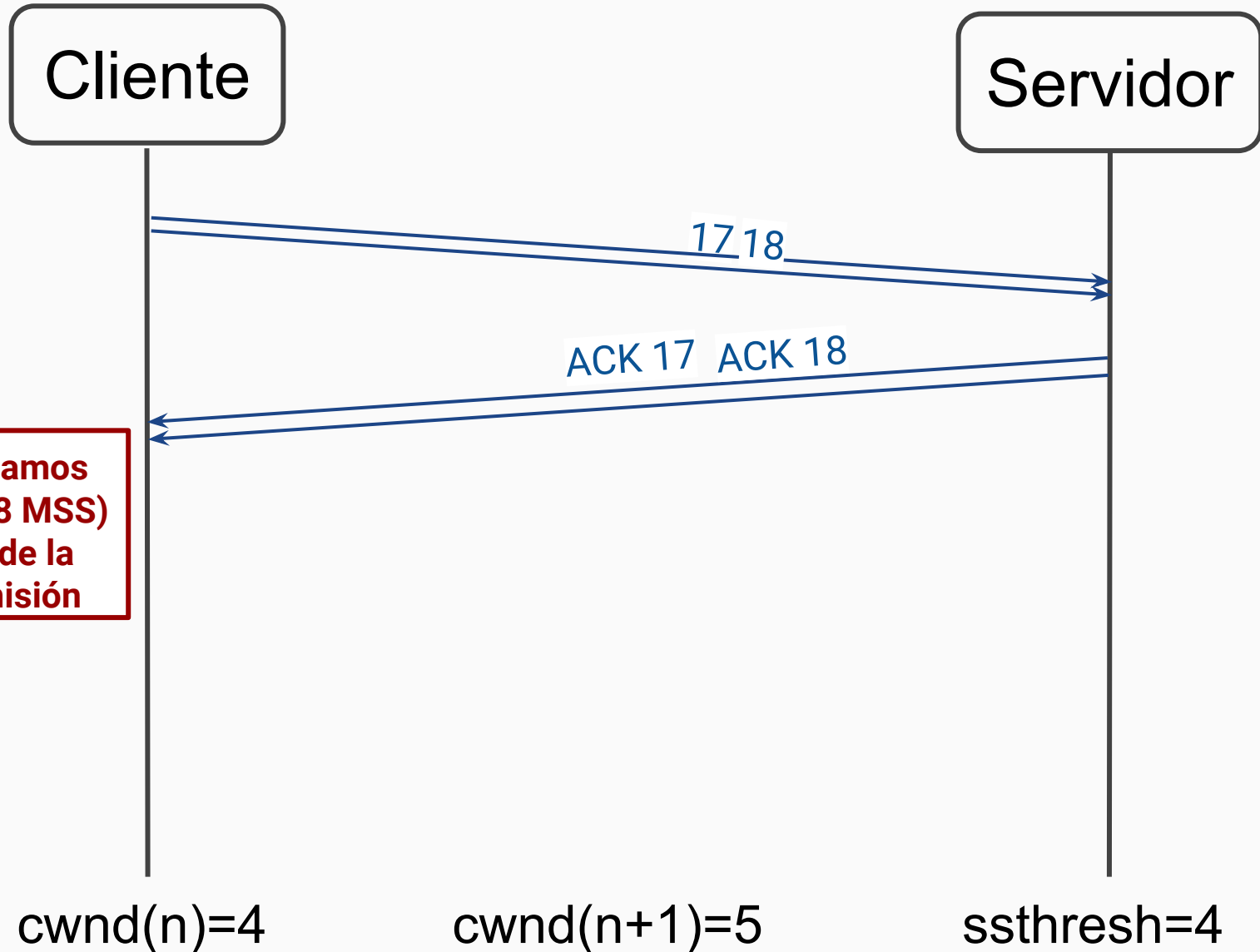
Slow start



Slow start



Slow start



Pérdida de paquetes

- **Ahora, veamos el caso de Reno**
- Supongamos que se pierde **UN SOLO** segmento de la ráfaga en cuestión.
- Y además, que estamos usando el algoritmo de control de congestión de Reno

TCP : control de congestión

Algoritmo TCP Reno

- Cuando se reciben 4 ACKs iguales, se retransmite el siguiente segmento
- $cwnd(n+1) = cwnd(n) / 2$
- $ssthresh = cwnd(n) / 2$
- Se continúa con la fase de CA en vez de pasar a **Slow Start**
- Esto se denomina **Fast Recovery**

TCP : control de congestión

Otros algoritmos

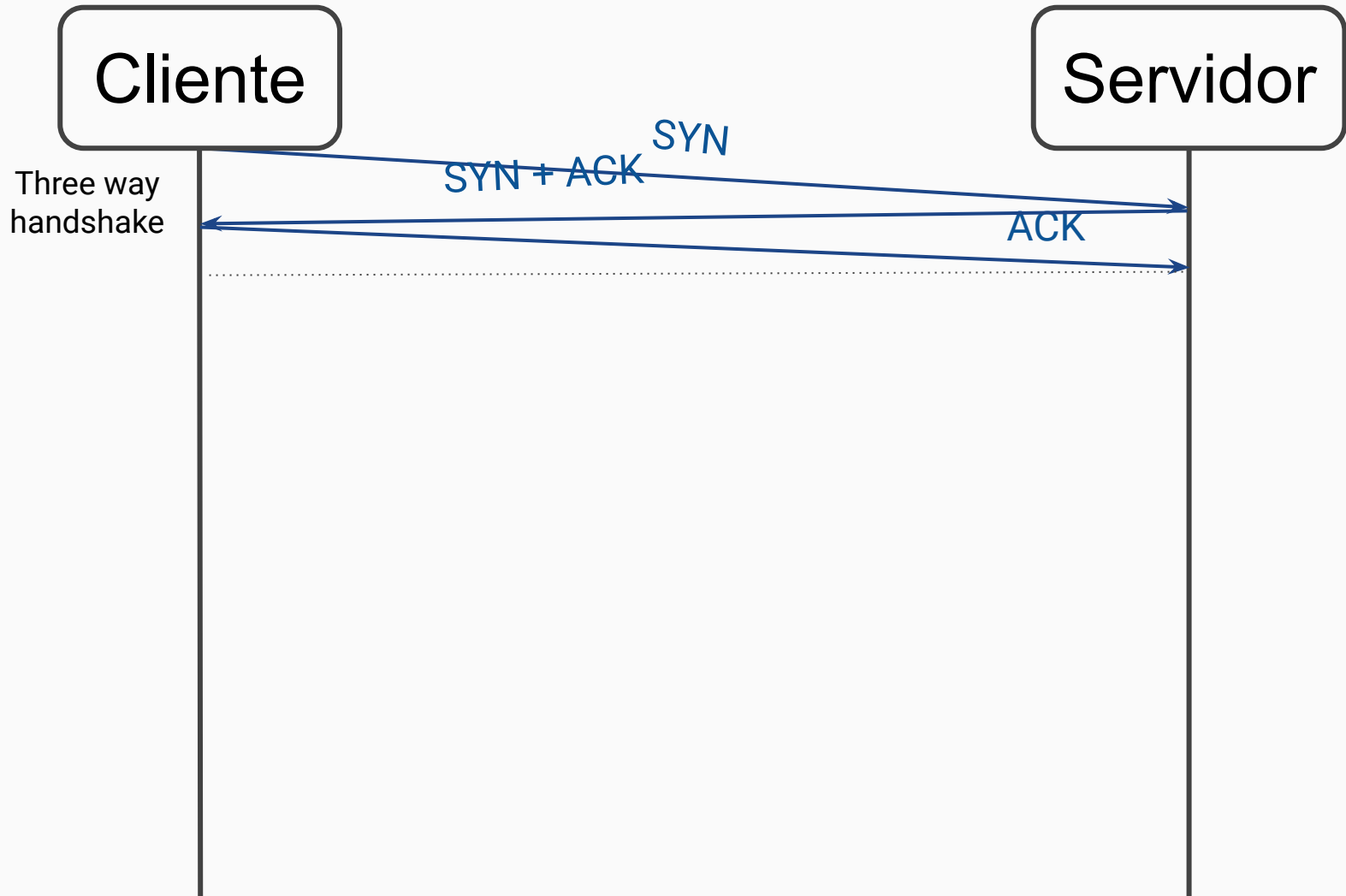
- TCP Tahoe y Reno son los más simples a nivel didáctico
- En la práctica ya no se utilizan. Hoy en día se usan algoritmos más modernos y complicados

TCP : control de congestión

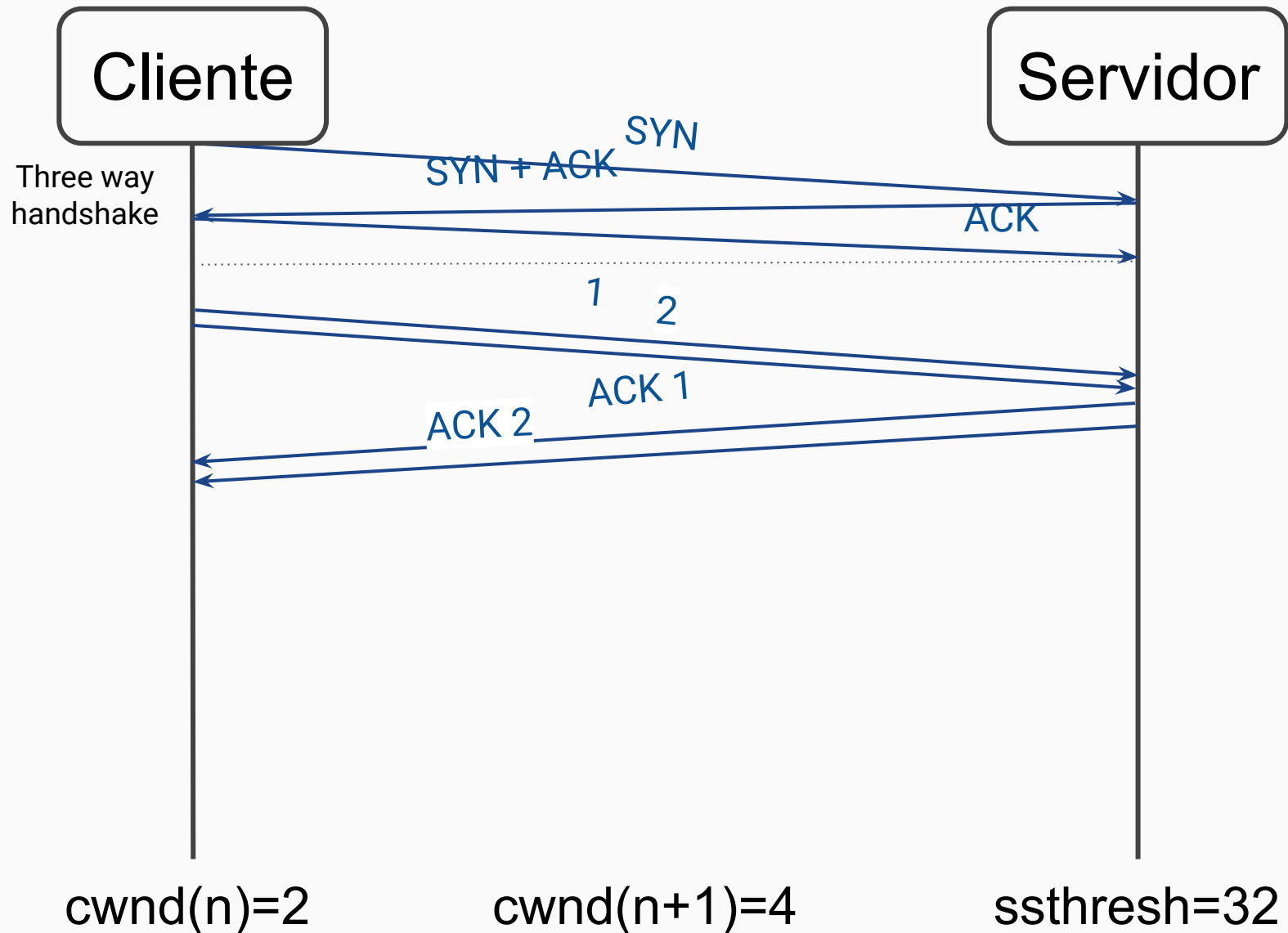
Otros algoritmos

- Algunos ejemplos
 - TCP Vegas
 - TCP New Reno (viene por default en FreeBSD)
 - TCP CUBIC (viene por default en Linux)

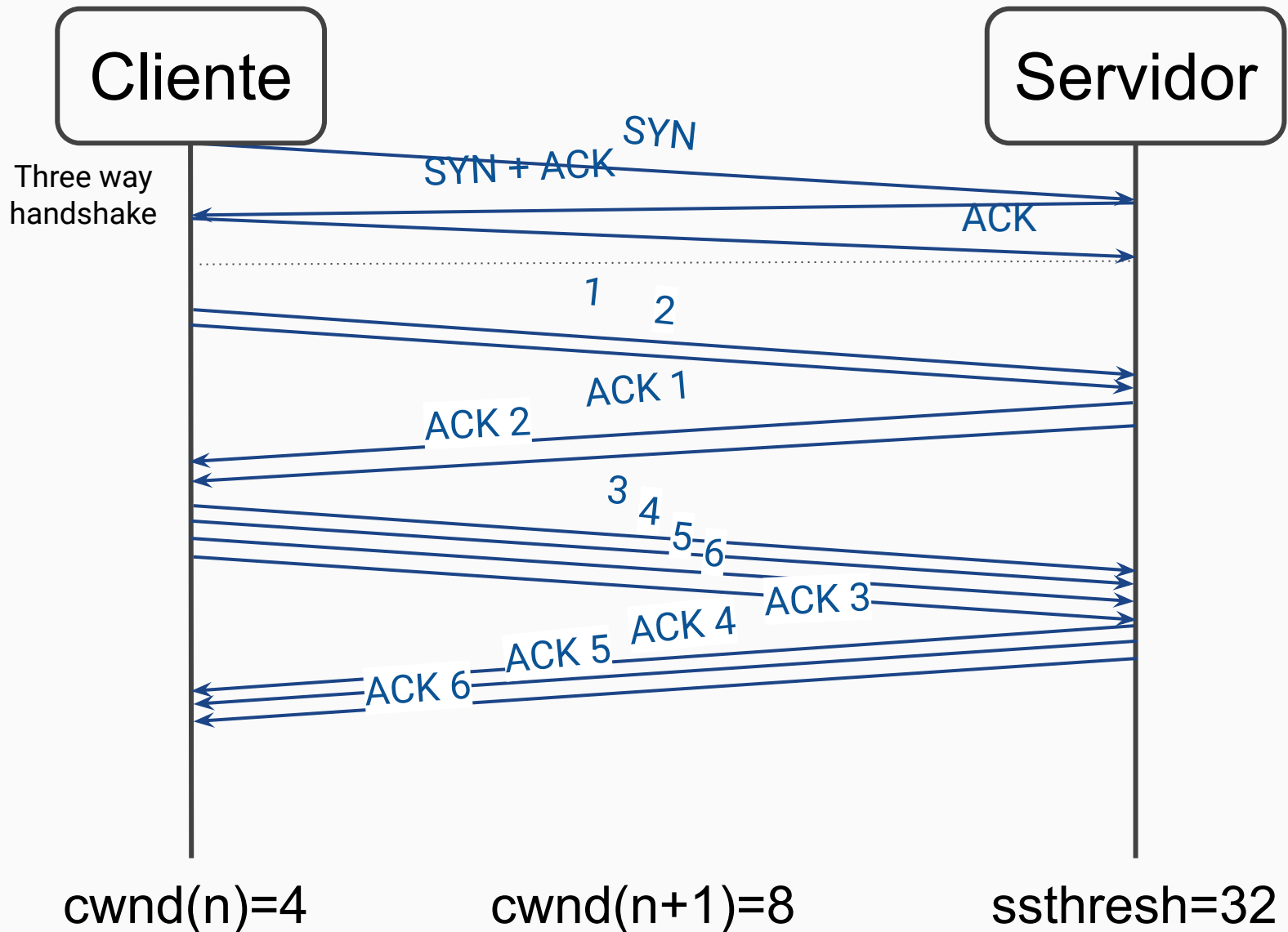
Three way handshake



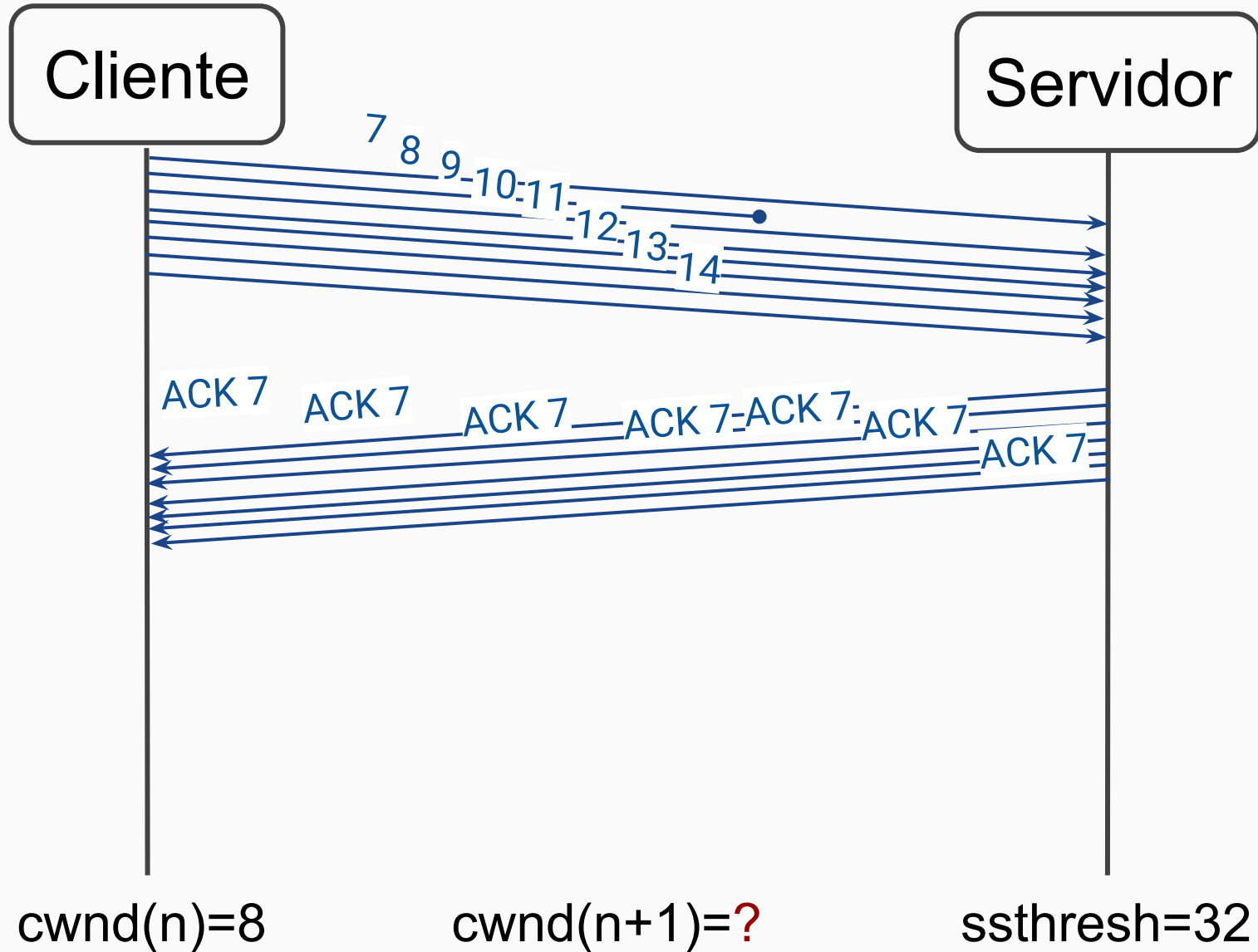
Slow start



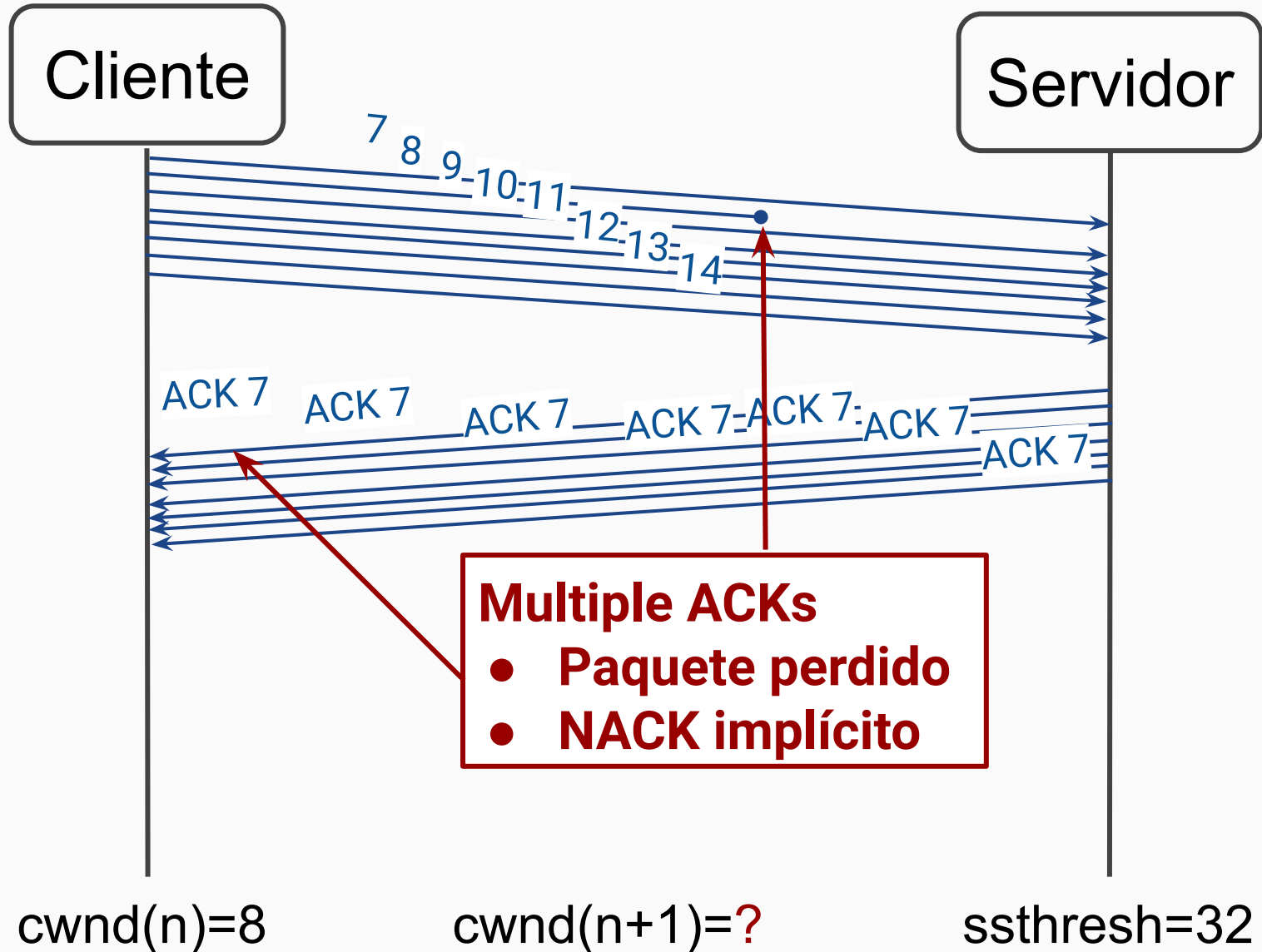
Slow start



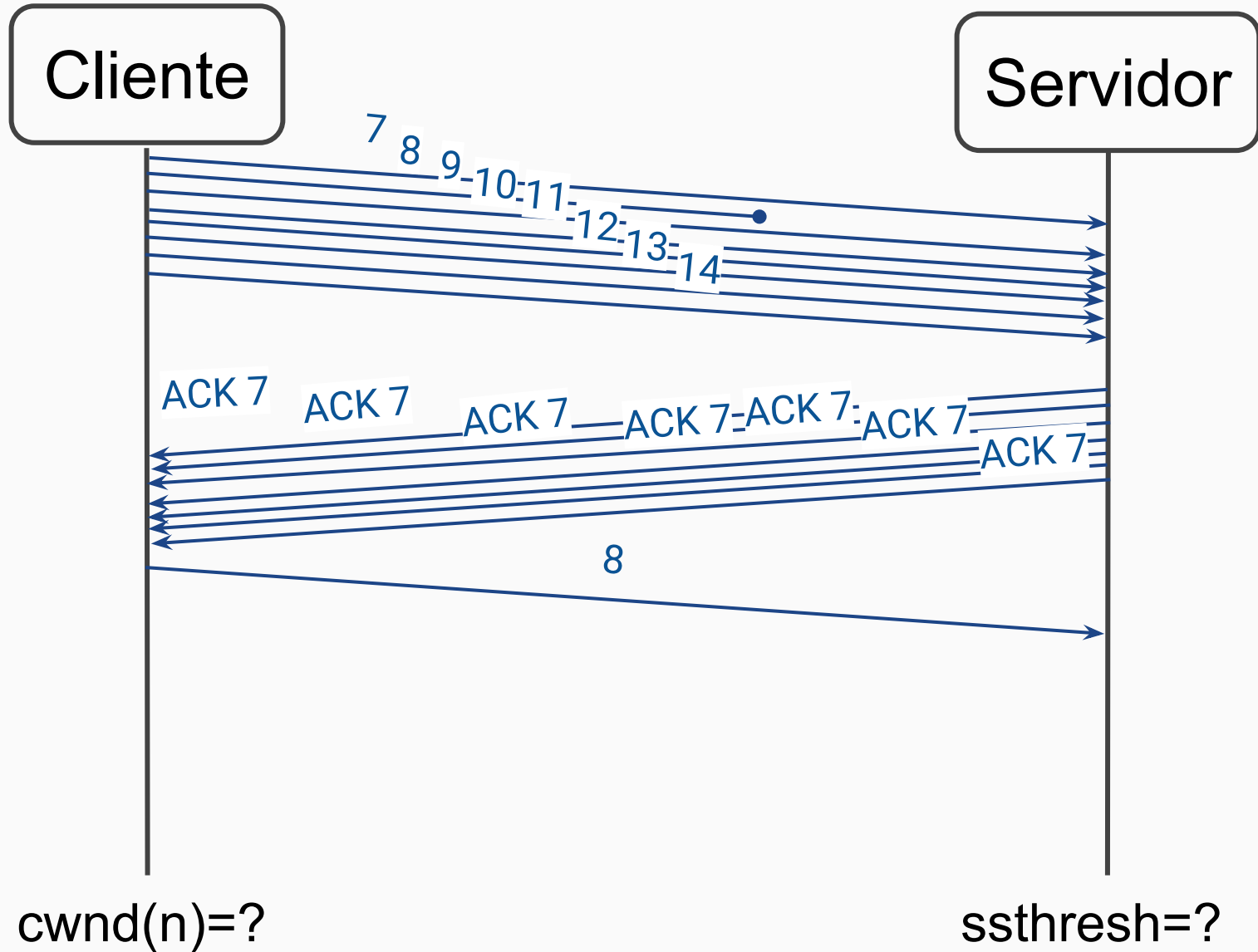
Slow start



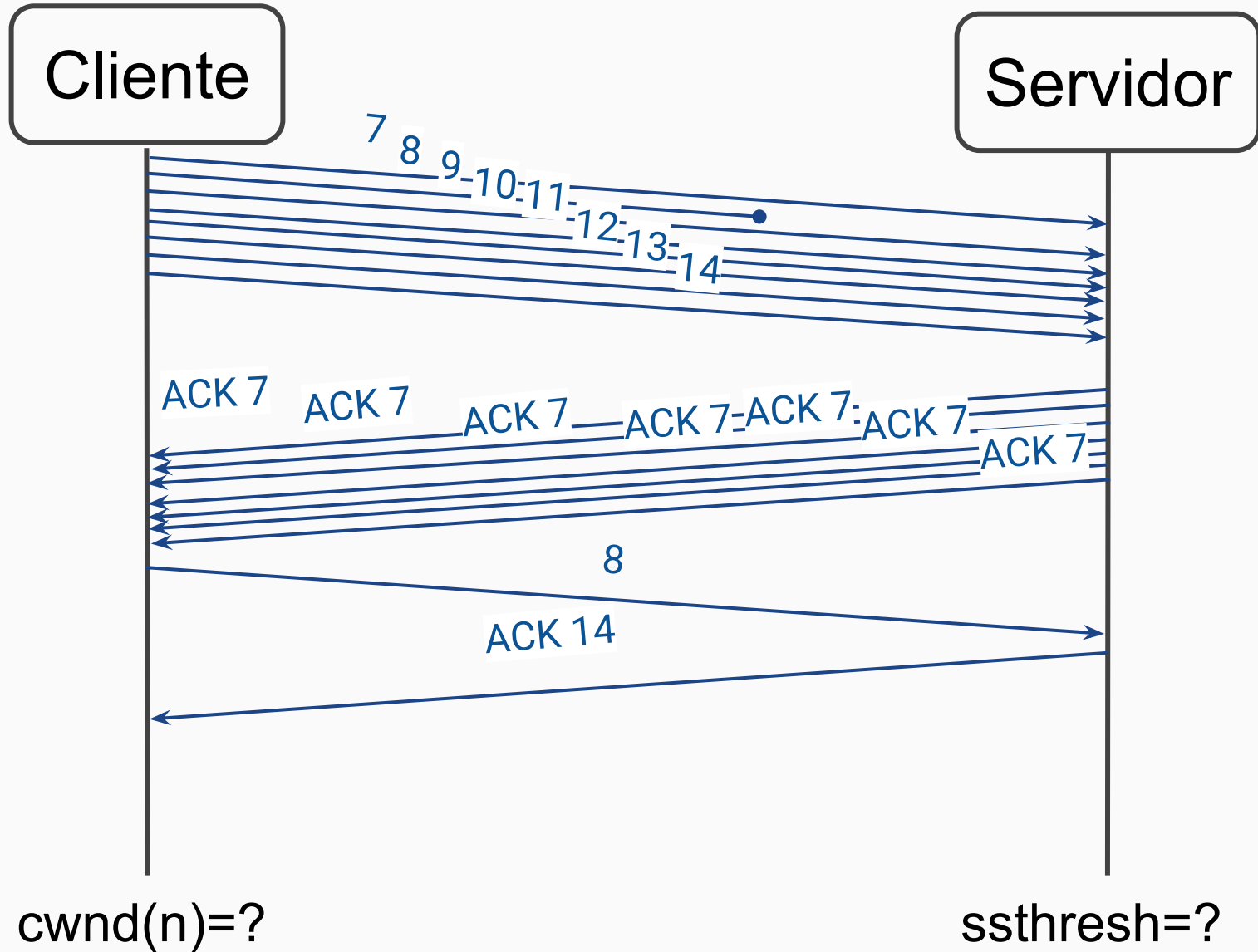
Slow start



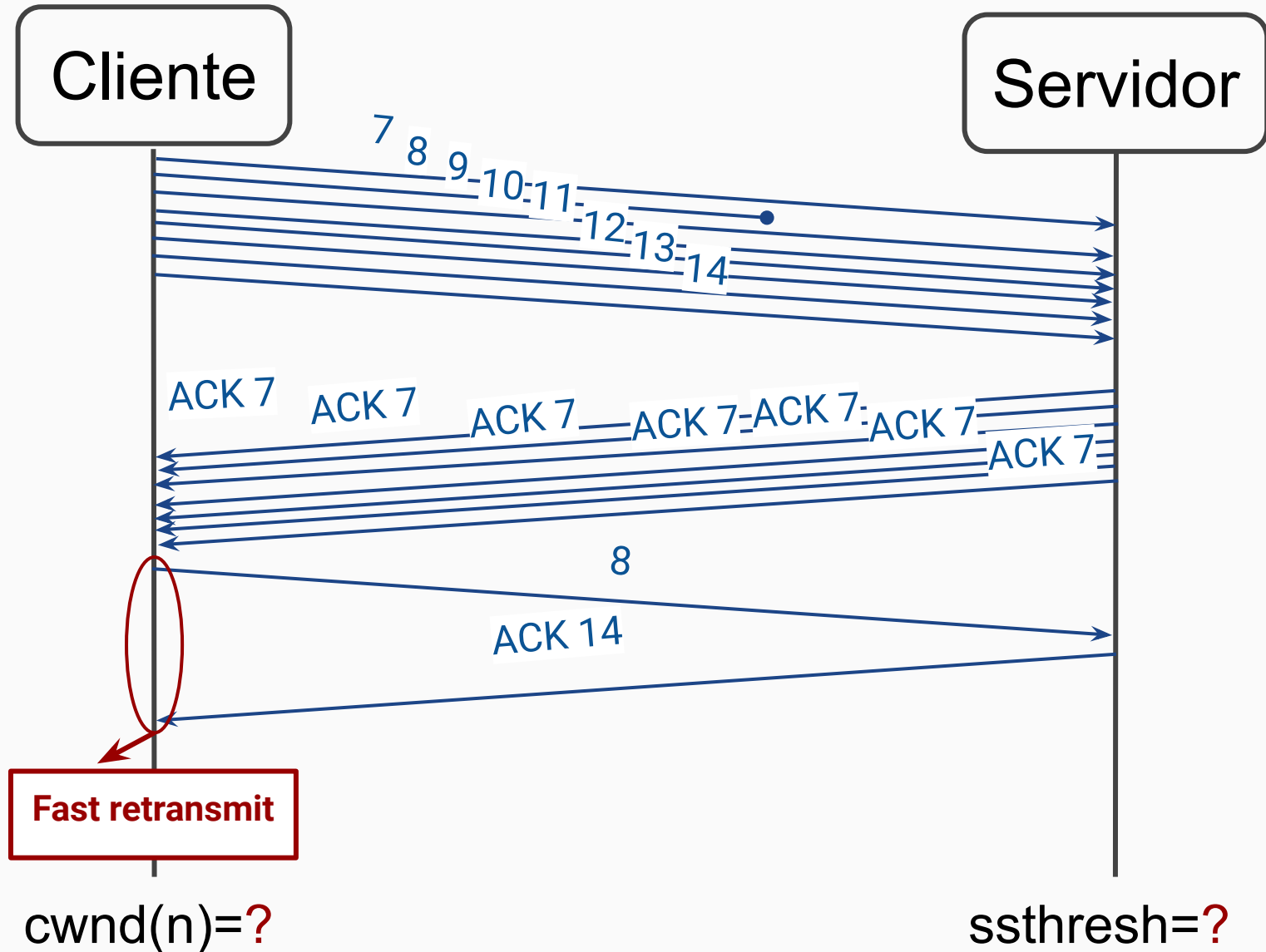
Slow start



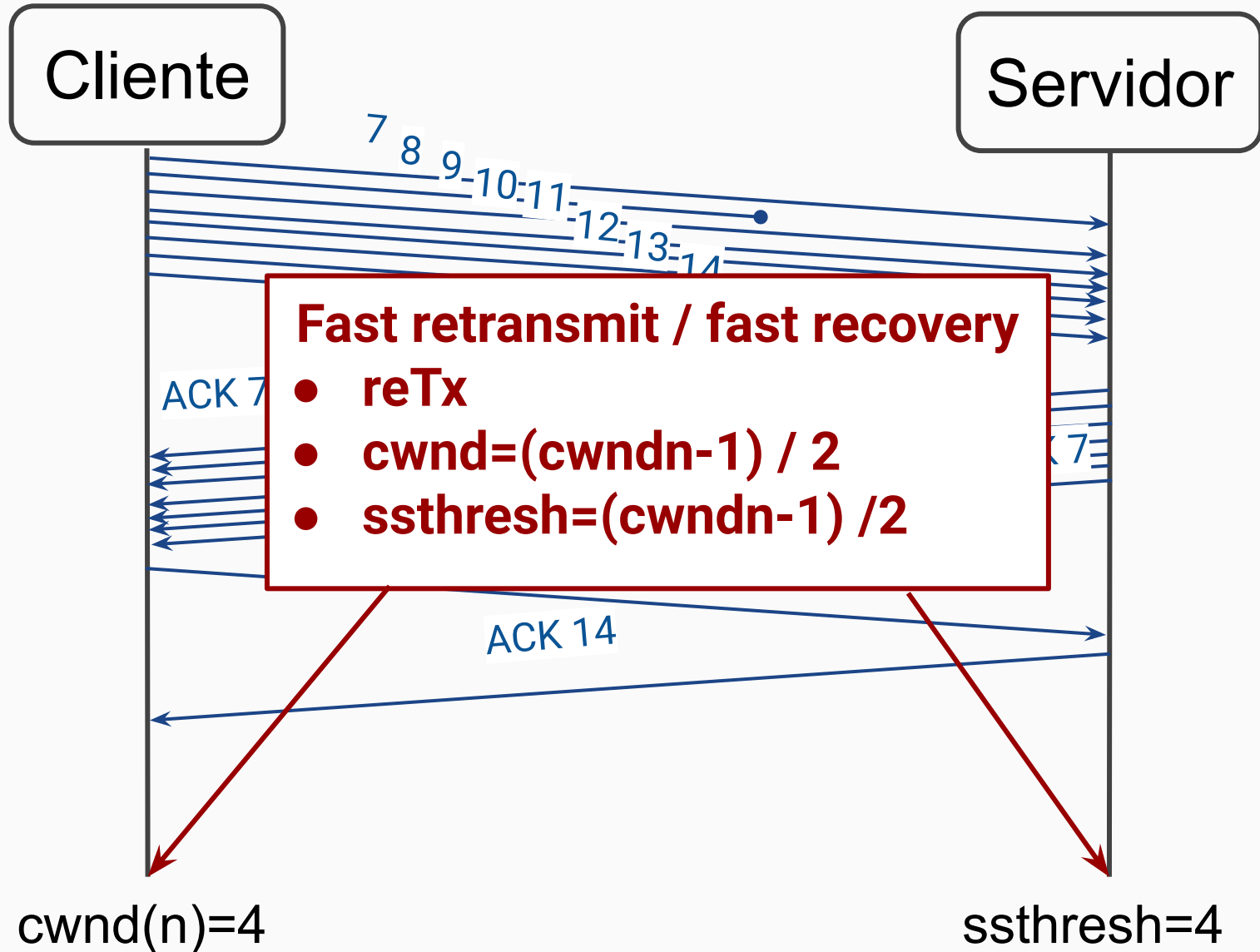
Slow start



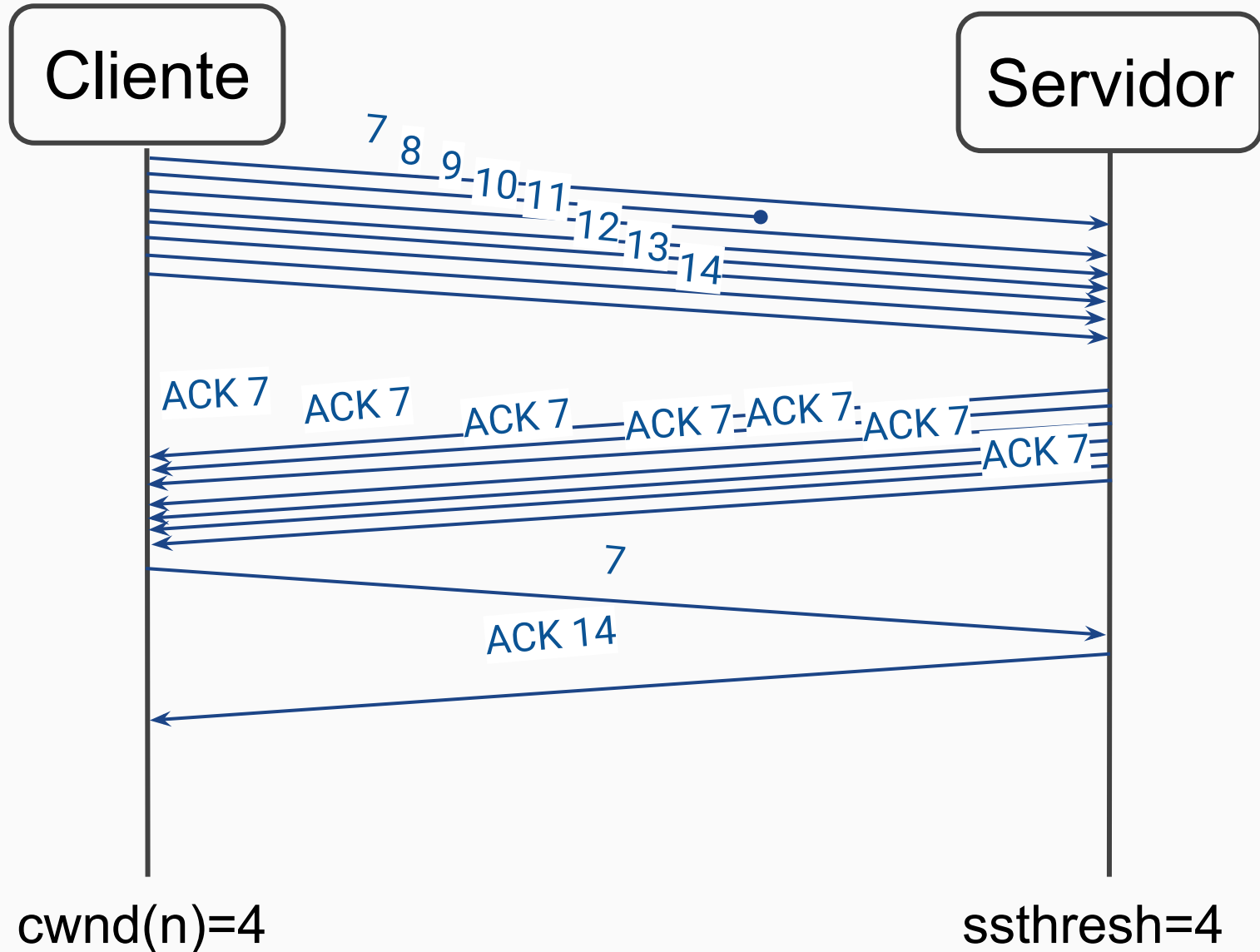
Fast retransmit



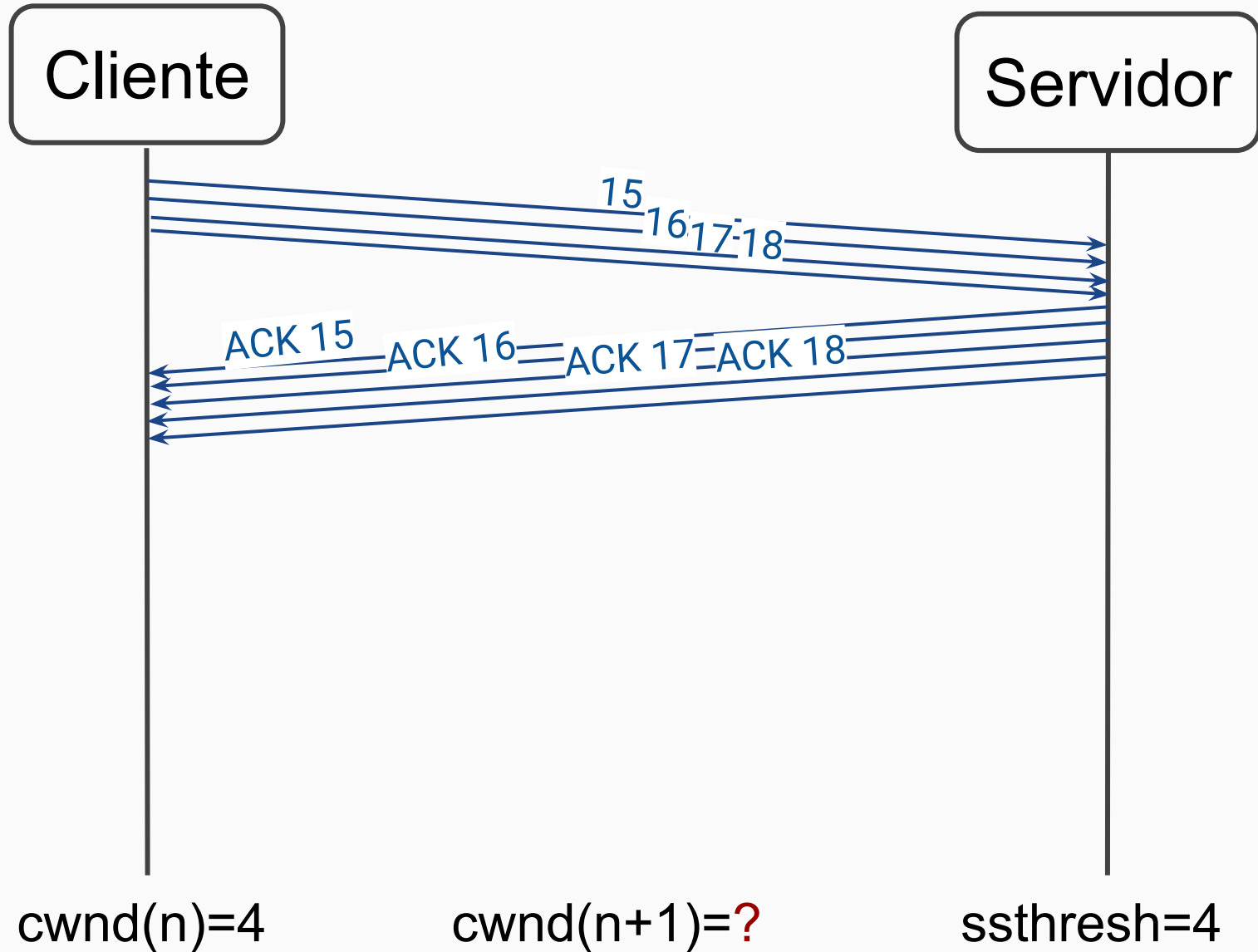
Fast retransmit



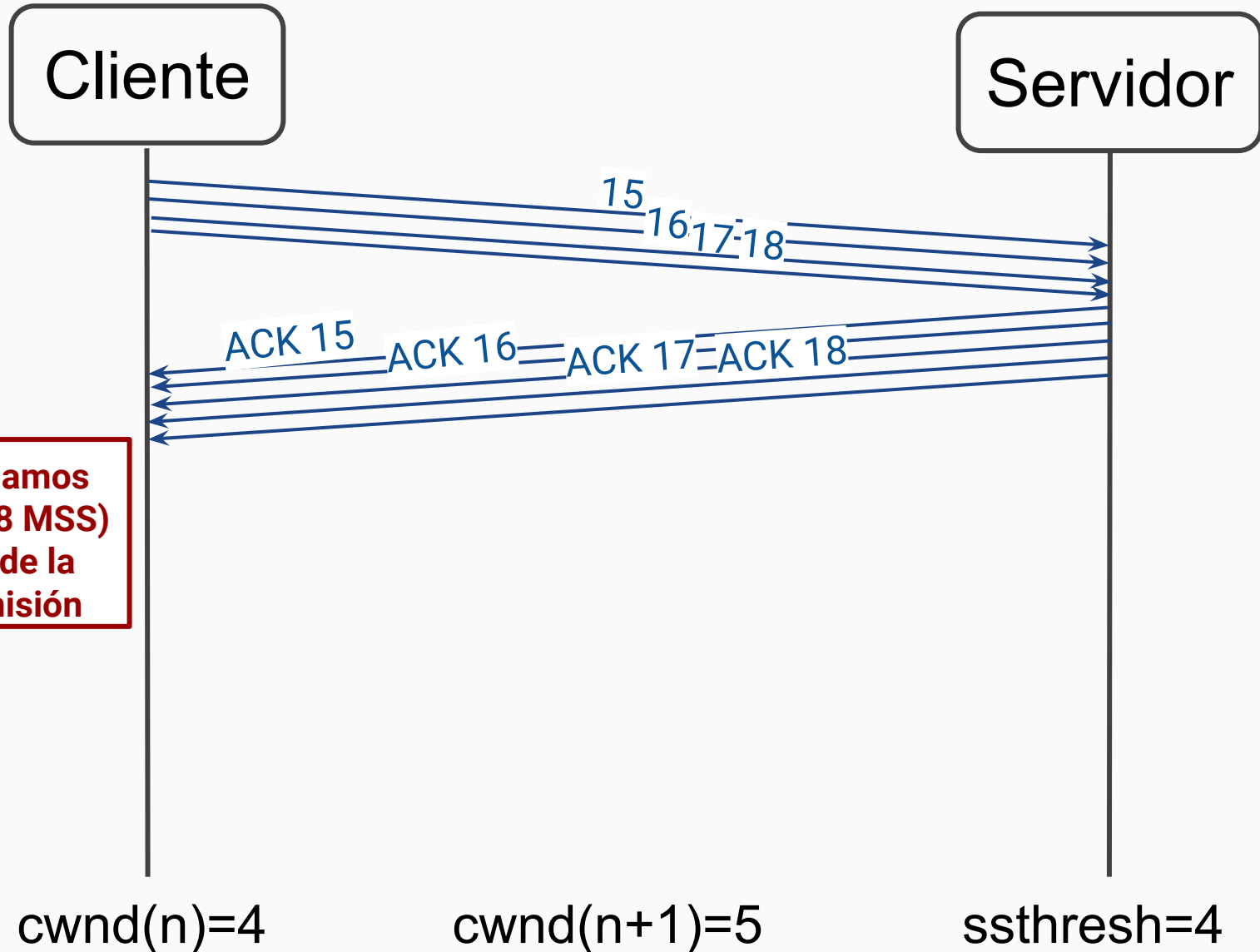
Fast retransmit



Fast retransmit



Fast retransmit



¿Preguntas?

Close TCP

TCP states

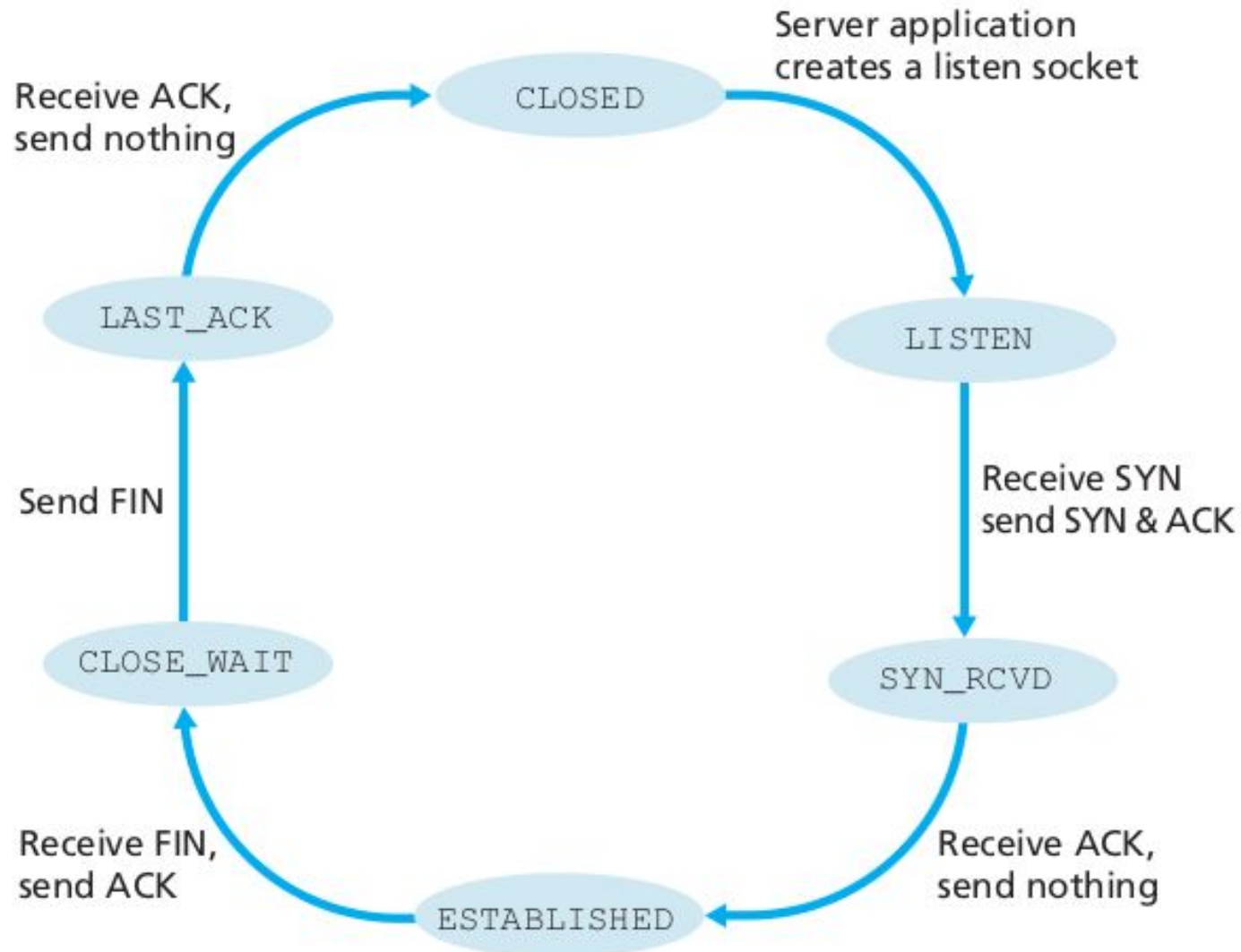


Figure 3.42 ♦ A typical sequence of TCP states visited by a server-side TCP

TCP states

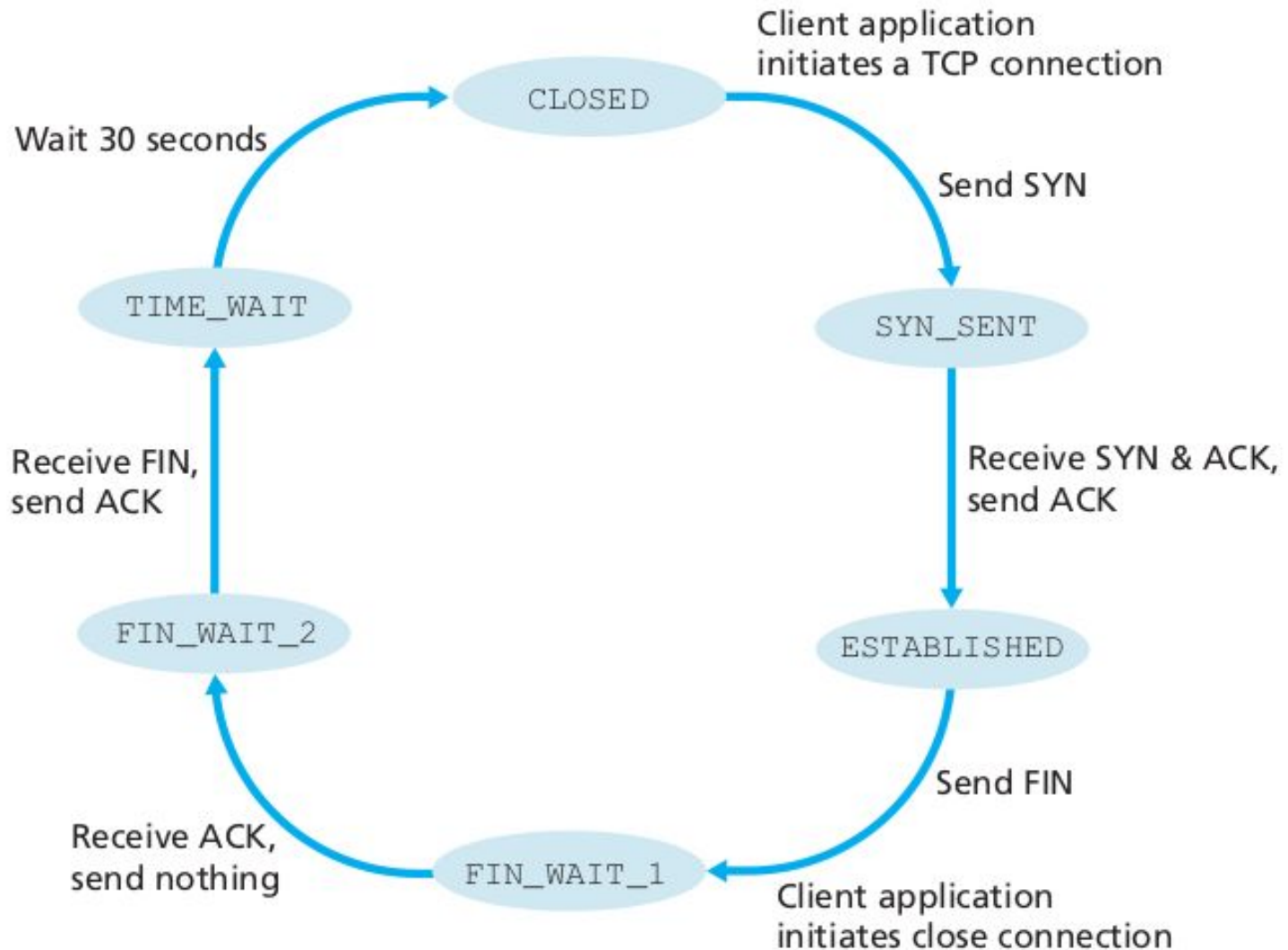
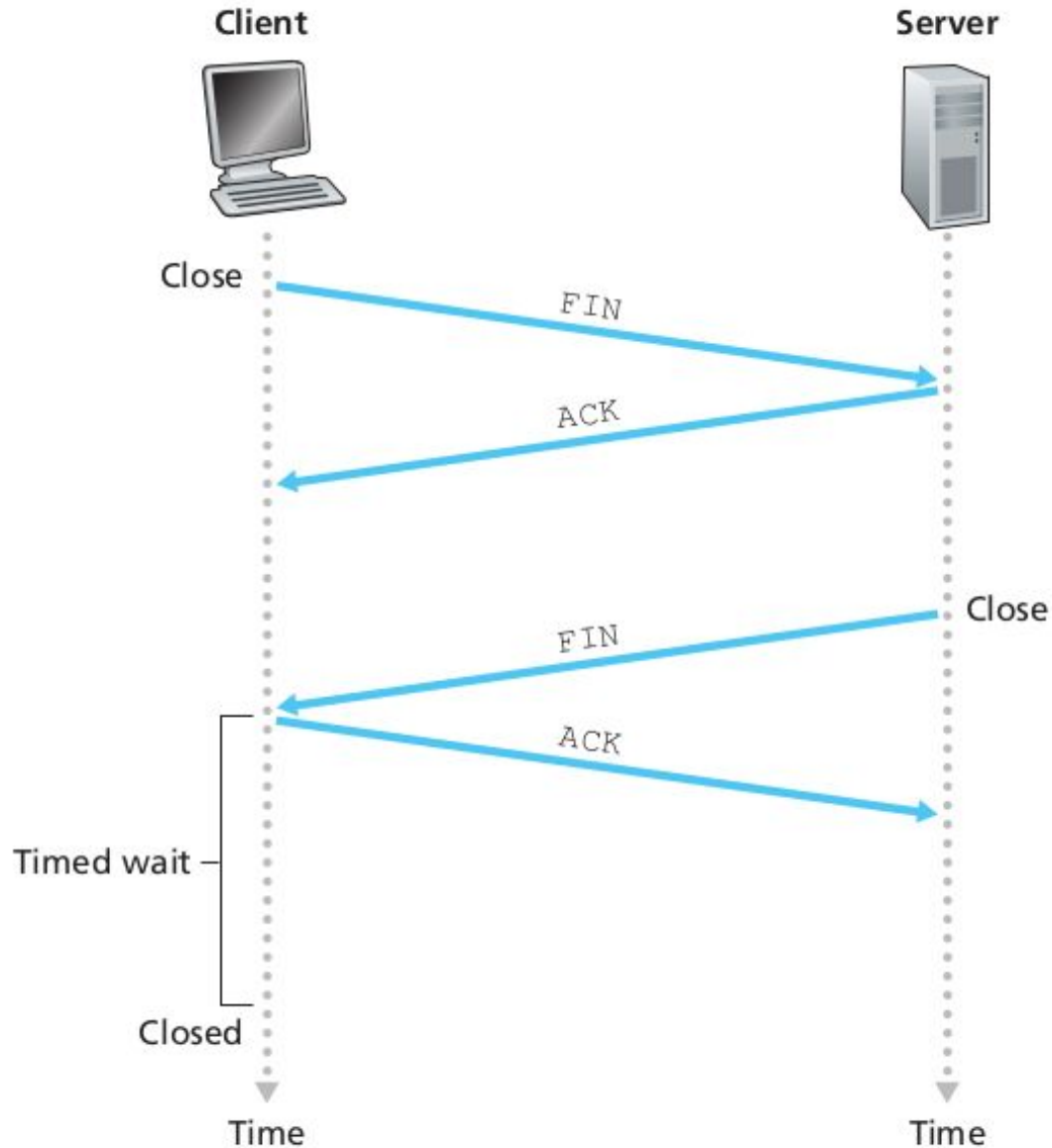


Figure 3.41 ♦ A typical sequence of TCP states visited by a client TCP

Close messages



Split TCP

Split TCP

- **Problema**
 - Latencia alta.

Split TCP

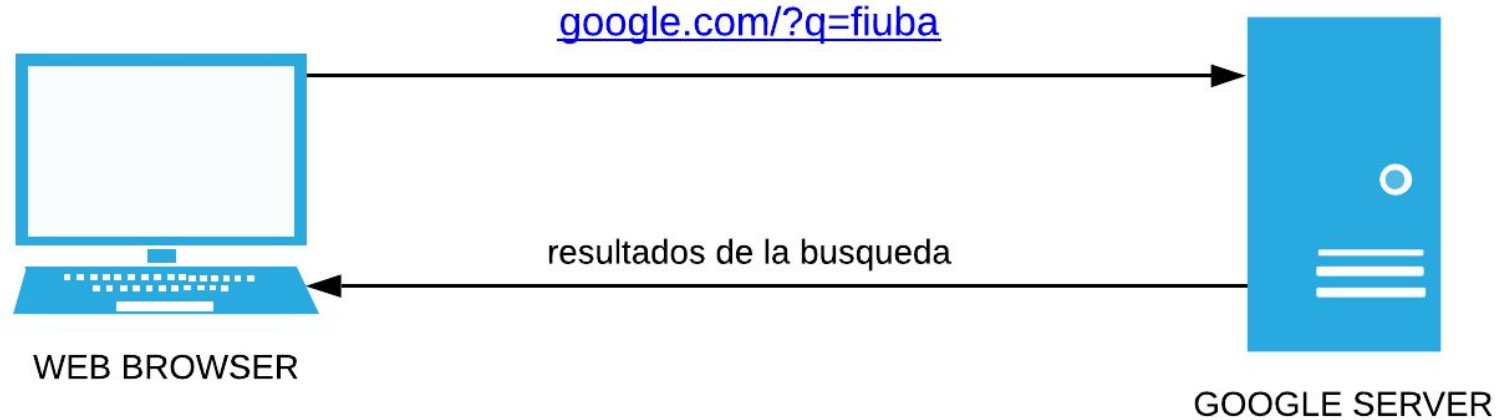
- **Problema**

- Latencia alta.
- Recursos dinamicos.

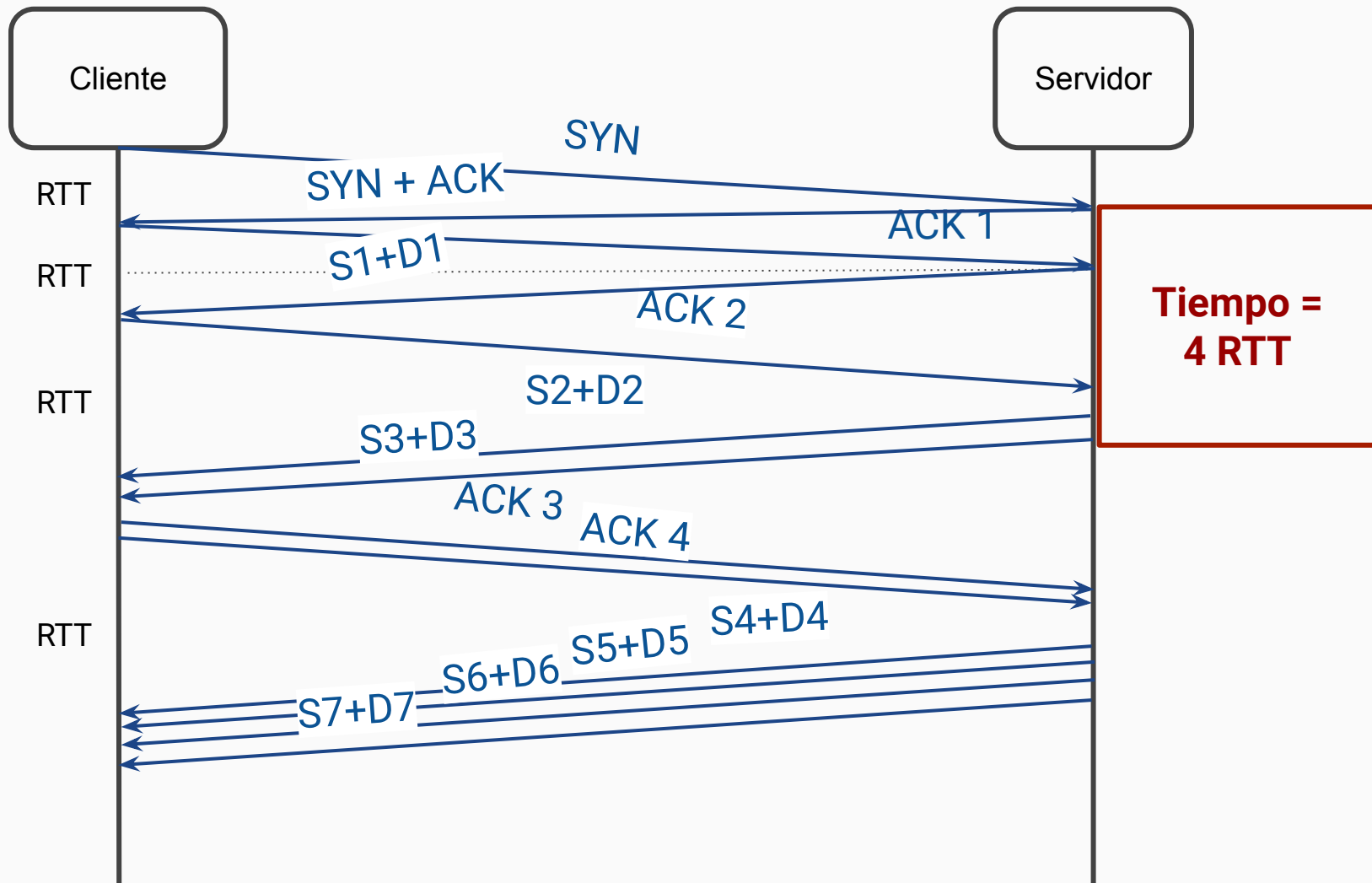
Split TCP - Ejemplo

- Ejemplo Kurose: HTTP search query
- “Typically, the server requires three TCP windows during slow start to deliver the response [Pathak 2010]”
- A estos 3 RTT tenemos que sumar el RTT para el three way handshake

Split TCP - Ejemplo



Split TCP - Ejemplo



Split TCP - Ejemplo

- Tardamos 4 RTTs en enviar el resultado de la query
- Se podría mejorar?

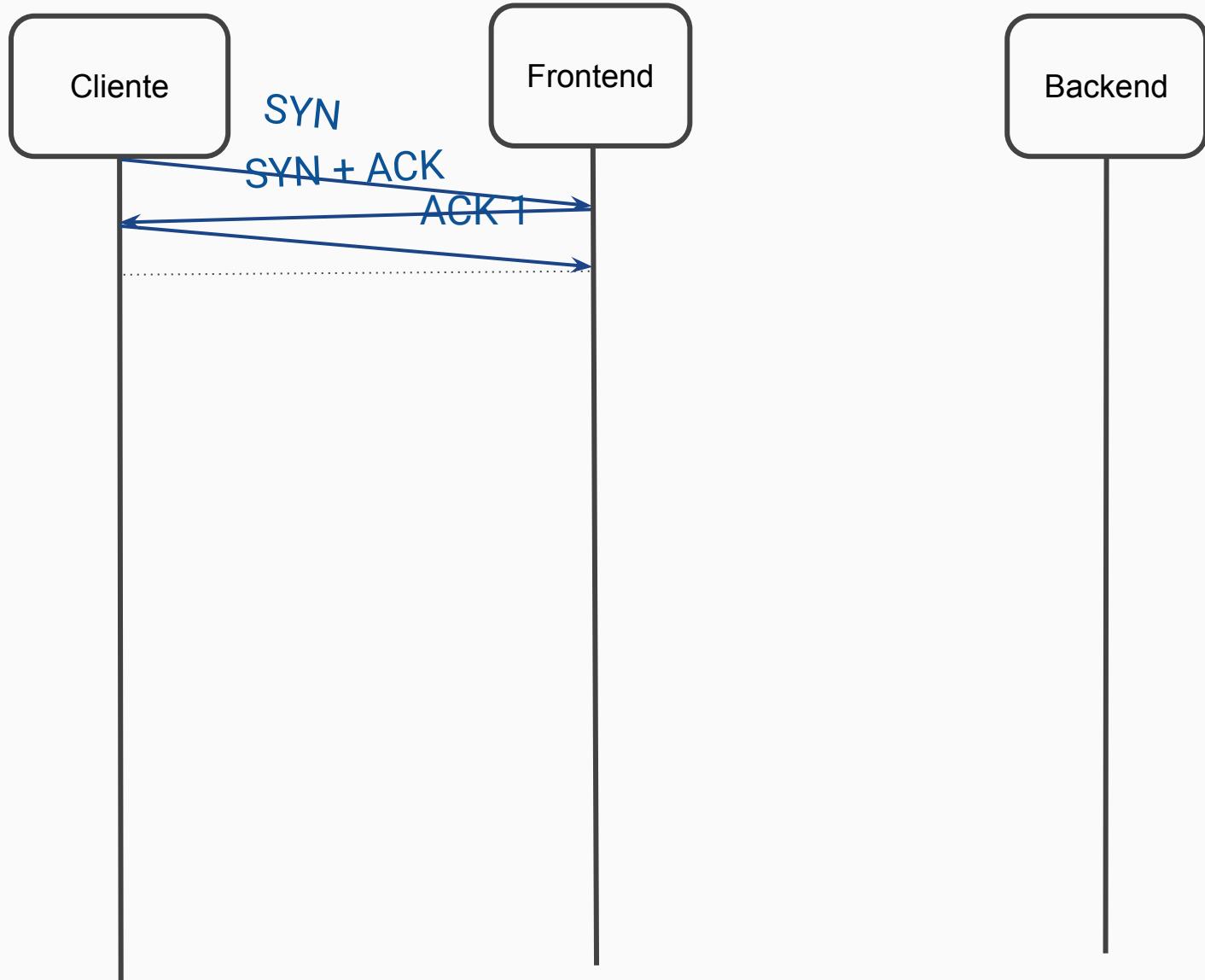
Split TCP

- Split TCP:
 - **Frontend:** cercano al cliente, ventanas más chicas

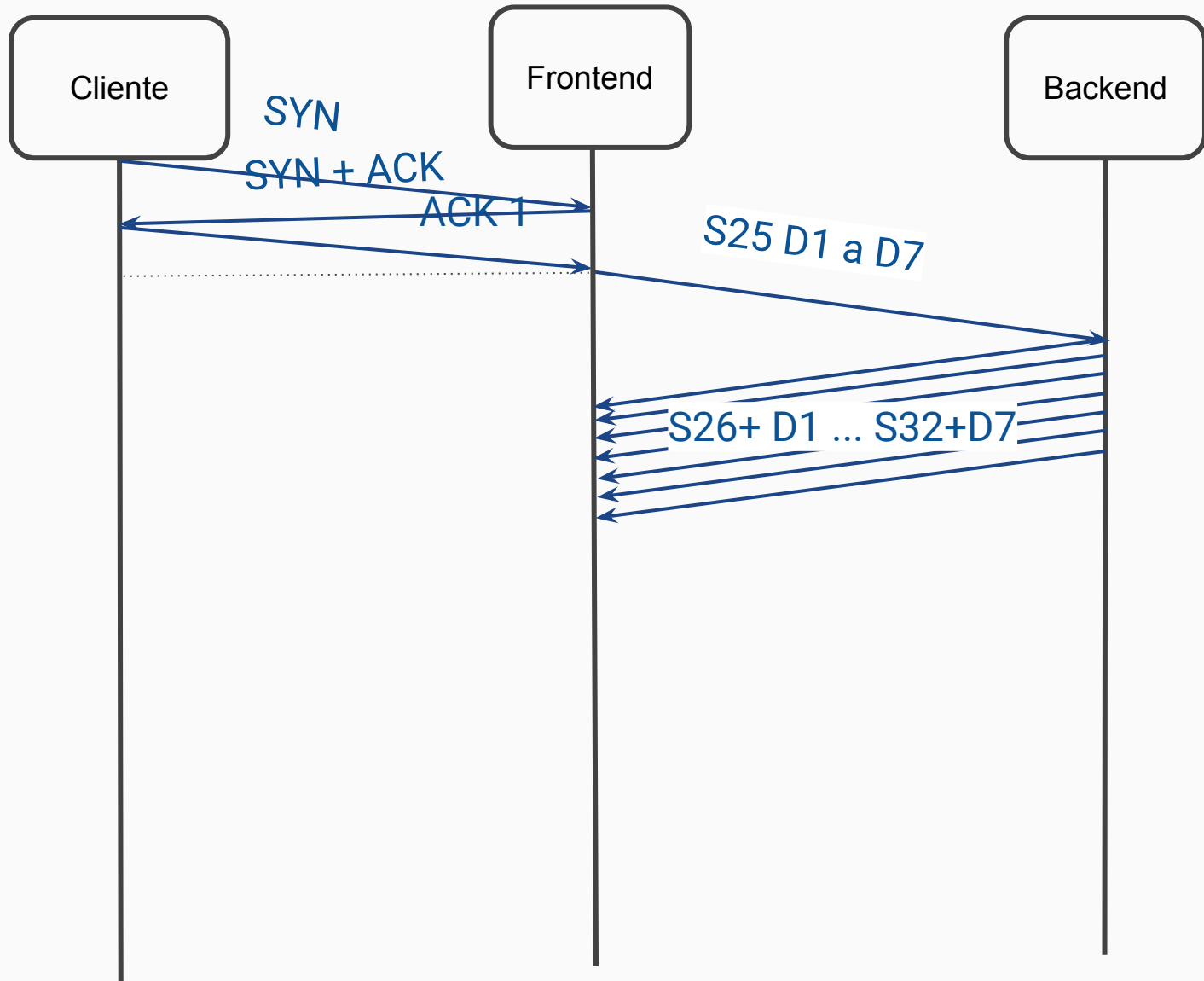
Split TCP

- Split TCP:
 - **Backend:** lejano al cliente
 - Tiene todo el tiempo una conexión activa con el **FE**
 - Ventanas más grandes

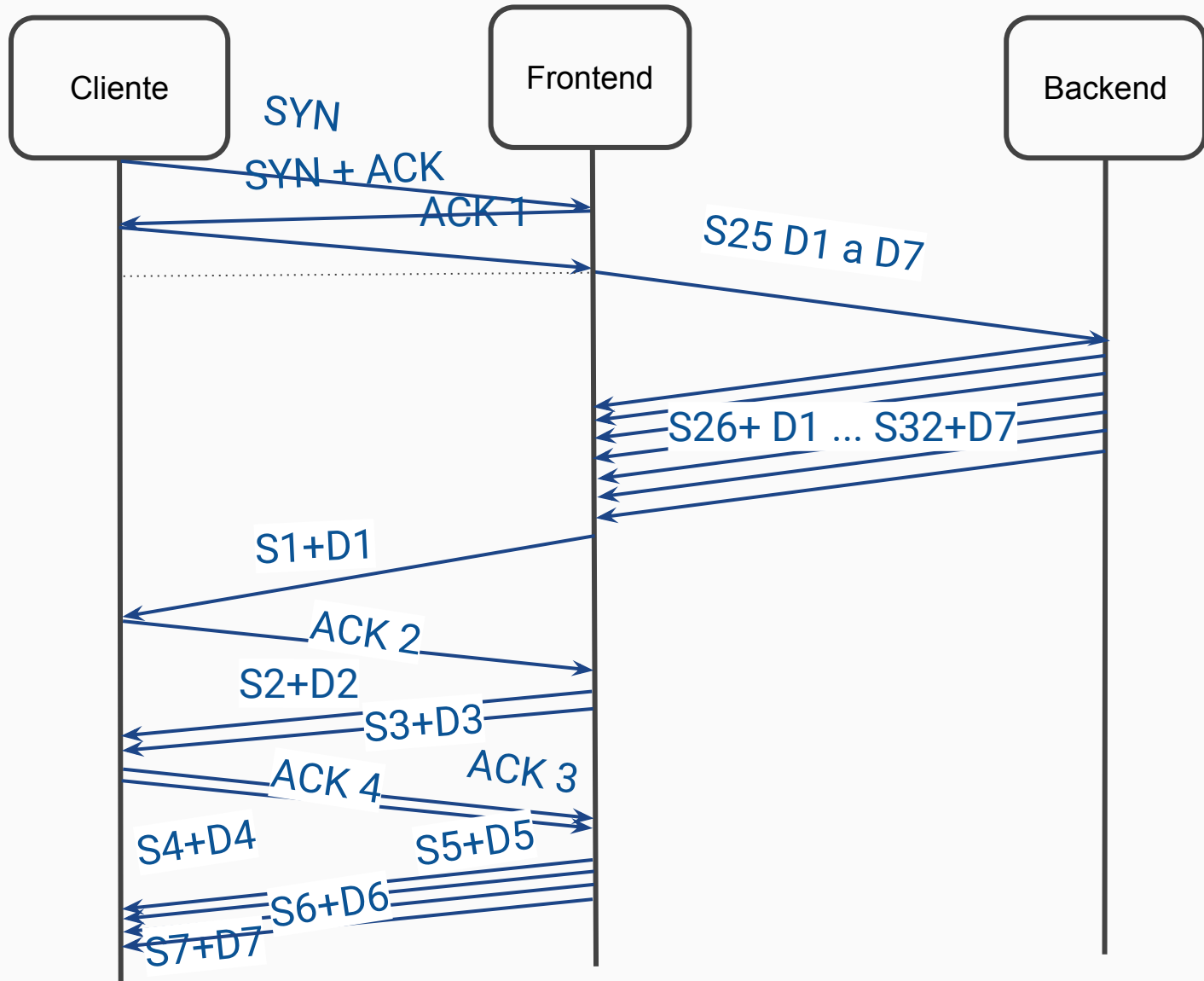
Split TCP - Ejemplo



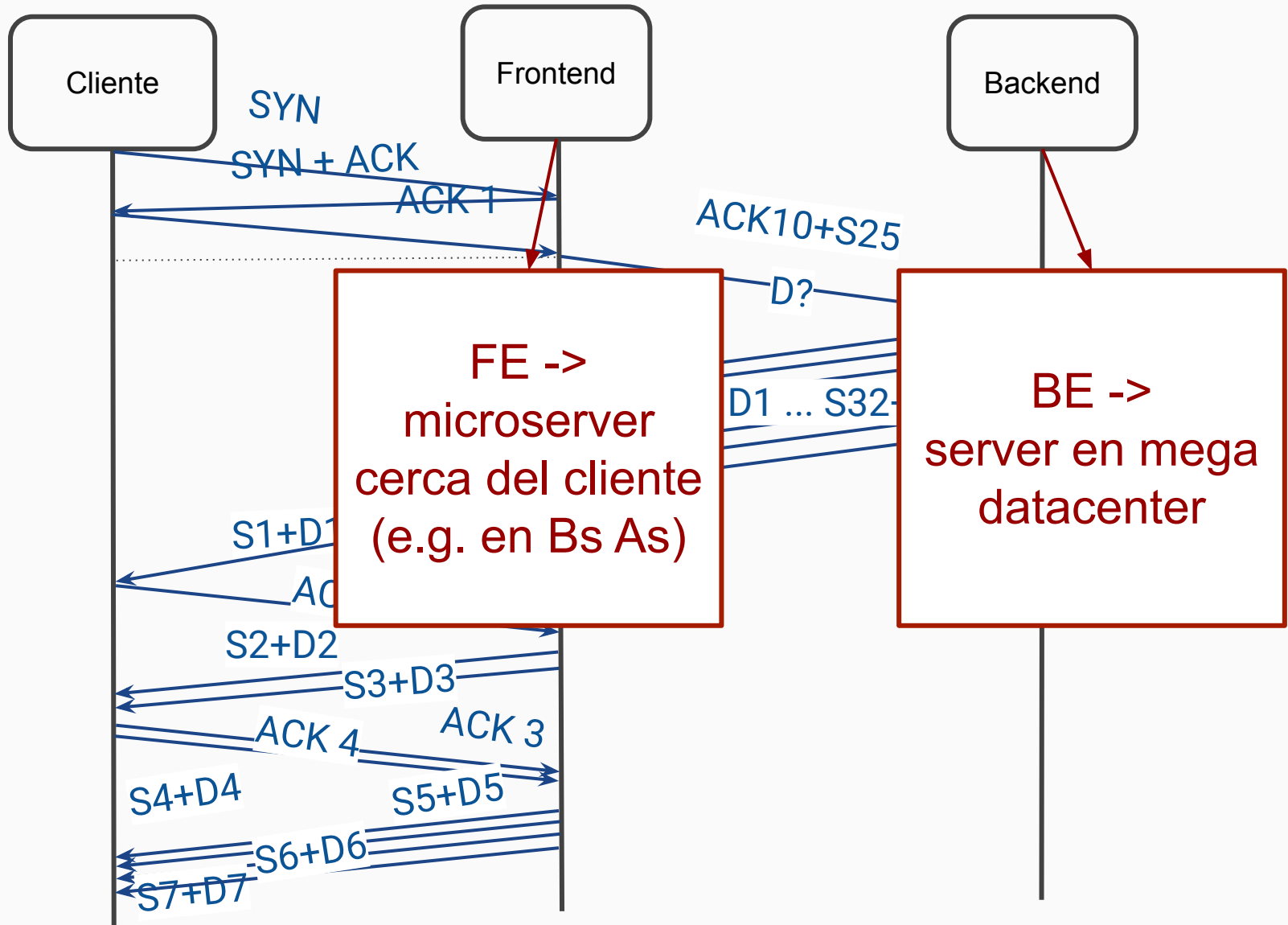
Split TCP - Ejemplo



Split TCP - Ejemplo



Split TCP - Ejemplo



Split TCP - Ejemplo

