

UNIVERSITA' DEGLI STUDI DI
NAPOLI FEDERICO II

Scuola Politecnica e delle Scienze di Base
Corso di Laurea in Ingegneria Informatica



elaborato di laurea

Il controllo di congestione BBR ed il suo impiego nella rete B4 di Google

Anno Accademico 2016/2017

candidato
Riccardo Carbone
matr. N46002406

Indice

1	Prefazione	2
2	Introduzione	3
3	I problemi degli approcci loss-based	4
3.1	Il limite del delivery rate	4
3.2	Il sending rate e la congestione di rete	5
3.3	Gli approcci loss based	7
4	Il controllo di congestione BBR	8
4.1	Il network path model	9
4.1.1	Bottleneck bandwidth	9
4.1.2	Round Trip propagation delay	13
4.2	I parametri di controllo	13
4.2.1	Pacing rate	14
4.2.2	Send quantum	14
4.2.3	Congestion window	15
4.3	Visione complessiva dell'algoritmo	16
4.3.1	Startup	17
4.3.2	Drain	18
4.3.3	ProbeBW	18
4.3.4	ProbeRTT	20
4.3.5	Evoluzione di un flusso BBR	21
5	Casi reali d'impiego del BBR	22
5.1	Cubic vs BBR	22
5.2	Fairness tra flussi BBR	24
5.3	Fairness tra flussi loss-based e BBR	25
5.4	Rete Google B4	25
6	Appendice A	28
6.1	Motivazione	28
6.2	Bic-TCP	28
6.3	Cubic	30

Prefazione

La scelta dell'argomento trattato all'interno dell'elaborato (il controllo di congestione BBR) è ricaduta su una delle tematiche che sono state affrontate durante il progetto #WikiTIM, svoltosi nel periodo di Giugno-Luglio 2017 nell'ambito di una collaborazione tra l'Università di Napoli Federico II, Wikimedia Italia e TIM.

Lo scopo di tale progetto era, ed è, quello di avvicinare il mondo dell'università alla più grande enciclopedia libera del Web: Wikipedia. E' infatti fondamentale per una maggiore valorizzazione del patrimonio informativo che essa contiene, far sì che ogni sua voce non sia solo chiara, ma corredata di riferimenti bibliografici a fonti attendibili ed accessibili a tutti.

Il progetto #WikiTIM ha avuto come risultato finale la valorizzazione delle voci relative al: controllo di congestione, ed alla gestione di rete.

In particolare, il gruppo di lavoro dell'Università Federico II di Napoli ha contribuito alla revisione delle voci:

- **Controllo di congestione:**
https://it.wikipedia.org/wiki/Controllo_della_congestione
- **Controllo di congestione in TCP:**
https://it.wikipedia.org/wiki/Controllo_della_congestione_in_TCP
- **Gestione di rete:**
https://it.wikipedia.org/wiki/Gestione_di_rete
- **Simple Network Management Protocol:**
https://it.wikipedia.org/wiki/Simple_Network_Management_Protocol

Introduzione

Nonostante gli enormi passi avanti che sono stati compiuti nel campo delle tecnologie usate nelle reti di calcolatori, sono ancora molteplici i casi in cui Internet non trasmette le informazioni alle velocità sperate.

Eppure molto s'investe nella ricerca e sviluppo di un'infrastruttura efficiente per la trasmissione dei segnali, che garantisca elevati livelli di delivery rate end-to-end (in ogni normale condizione).

Come dunque si può immaginare, la causa principale non è nè l'uso di tecnologie inadeguate, nè la mancanza di conoscenza nel loro uso. Il reale freno alle potenziali velocità raggiungibili, nasce da una scelta effettuata negli anni '80, quando fu definito il primo schema per il controllo di congestione del protocollo TCP.

Tale scelta riguardava l'identificare il fenomeno di congestione (in senso stretto) con il fenomeno di perdita dei pacchetti.

Tuttavia, nonostante i due concetti siano molto differenti, al tempo tale scelta aveva senso solo in virtù delle limitazioni tecnologiche.

Oggi anche la migliore versione loss-based di TCP: Cubic (App. A), fa fatica a garantire un delivery rate elevato. Basti pensare che per poter sostenere 10Gbps con un RTT di 100ms, è necessario un tasso di perdita dello 0.000003%, e con lo stesso RTT ed un tasso di perdita dell' 1% sono sostenibili solo 3Mbps [RSc17, p. 12].

Lo scopo di quest'elaborato è dunque quello di presentare un'alternativa per tali approcci, che non sia più basata sull'evento di packet loss, ma su una diversa e precoce rilevazione della congestione. L'alternativa che è stata approfondita è il nuovo controllo di congestione BBR.

Prima di iniziare, è lecito far presente che lo sviluppo di tale soluzione è ancora in corso, per cui alcuni aspetti di questa prima versione potrebbero subire delle modifiche nel corso del tempo.

I problemi degli approcci loss-based

In questo primo capitolo, illustreremo come mai un approccio guidato da eventi di packet loss, non riesce a garantire un buon delivery rate.

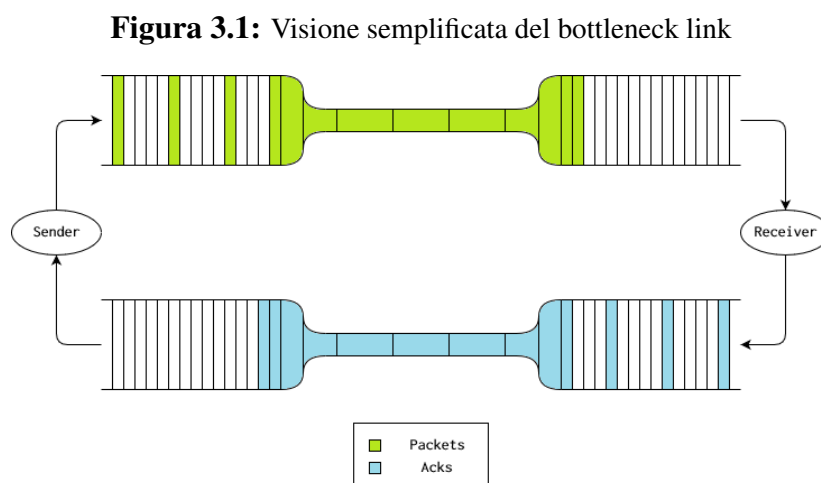
Il primo passo è comprendere le reali potenzialità del delivery rate.

3.1 Il limite del delivery rate

Nel momento in cui due end-point iniziano a comunicare attraverso la rete Internet, i loro pacchetti attraverseranno una serie di collegamenti (link) intermedi, ognuno dei quali sarà caratterizzato da una propria larghezza di banda disponibile.

In ogni istante, per ogni verso della comunicazione¹, esiste esattamente un solo bottleneck link: il collegamento con larghezza di banda disponibile più bassa. E' dunque necessario adeguarsi alle velocità di quest'ultimo, che pertanto determinerà il limite superiore del delivery rate end-to-end.

La seguente figura fornisce un'immagine semplificata ma coincisa, sul ruolo che riveste il bottleneck link nella comunicazione:



E' fisicamente impossibile per il delivery rate operare al di sopra del bottleneck rate.

¹In una rete a datagrammi come Internet non è detto che il bottleneck link del receiver sia lo stesso del sender. Dato che ogni pacchetto può in generale seguire un percorso diverso rispetto agli altri

3.2 Il sending rate e la congestione di rete

Dalla definizione del bottleneck rate ($BtlBw$) discende subito la seguente affermazione:

Un sending rate > bottleneck rate porta ad una crescita del bottleneck buffer, con conseguente aumento dei ritardi

Per poter caratterizzare il ritardo end-to-end introduciamo un altro parametro caratteristico del path di rete, ovvero il Round-trip-propagation-delay ($RTprop$).

L' $RTprop$ costituisce il minimo ritardo di propagazione, ovvero il limite inferiore dell'RTT per una data comunicazione.

A questo punto, si fa notare che l'influenza che ha il mittente nei confronti dello stato della rete (caratterizzato dallo stato dei buffer intermedi) dipende dal suo sending rate:

- Se il sending rate fosse minore o prossimo al bottleneck rate, il mittente introdurrebbe nella rete un amount inflight² tale da non contribuire alla crescita del bottleneck buffer. Per cui in tali condizioni l' $RTT = RTprop$:

Figura 3.2: Delivery rate and round trip time vs inflight part 1

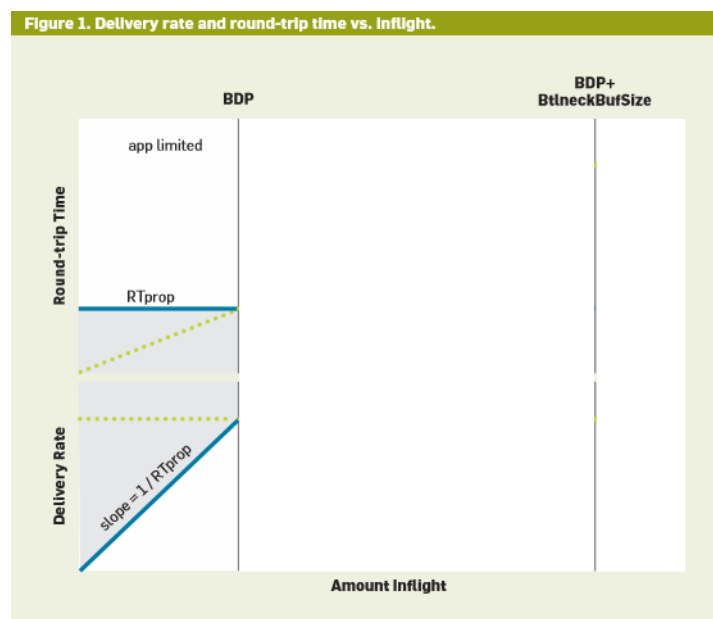


Figura tratta da [Car+17b, p. 60]

- Se il sending rate fosse superiore al bottleneck rate, il mittente introdurrebbe nella rete un amount inflight superiore alla capacità di quest'ultima. L'RTT aumenta per via della crescita del bottleneck buffer, e per bilanciare l'amount inflight così da fissare il delivery rate sul bottleneck rate:

²Quantità di dati inviata dal sender, ma non ancora riscontrata

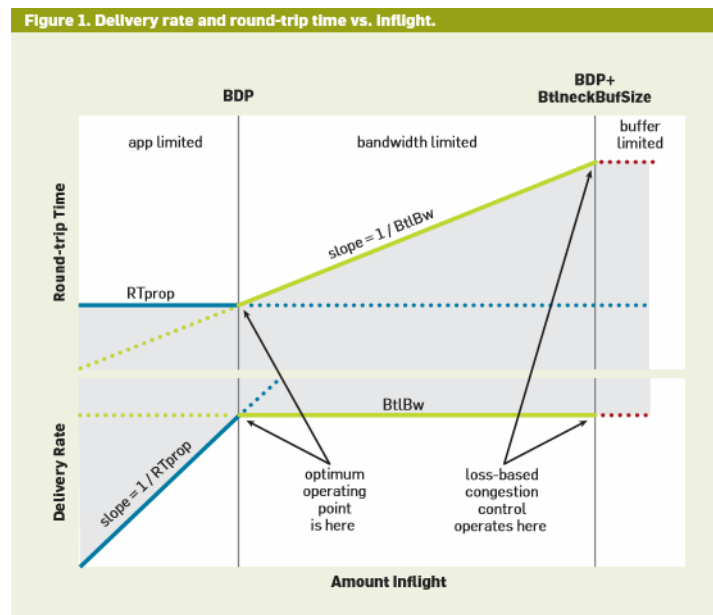
Figura 3.3: Delivery rate and round trip time vs inflight

Figura tratta da [Car+17b, p. 60]

Questi due casi sono molto importanti, in quanto nel primo il mittente non contribuisce all'aumento dei ritardi, mentre nel secondo sì. E' possibile quindi definire la condizione limite, che separa le due situazioni, introducendo il *BDP* (bandwidth delay product [Car+17b, p. 59]):

$$BDP = BtlBw * RTprop$$

Esso definisce la capacità di rete, ovvero il massimo amount inflight che il mittente può introdurre nella rete in un *RTprop* al massimo delivery rate, senza contribuire né alla generazione di code, né all'espansione dei ritardi.

Osservazione: interpretazione fisica dei parametri caratteristici

Per fornire una immagine più coincisa di tali parametri, possiamo immaginare una pipe virtuale tra il mittente ed il destinatario. Il *BtlBw* definisce il diametro della pipe, l'*RTprop* la sua lunghezza, mentre il *BDP* la sua capacità (per definizione).

Ora abbiamo tutti gli elementi per comprendere in che condizione il mittente contribuisce ad una congestione di rete, basta trasmettere in un *RTprop* una quantità di dati superiore al *BDP* corrente (pipe overflow).

3.3 Gli approcci loss based

Tali approcci limitano il sending rate del mittente, ponendo un limite superiore per l'inflight amount: la congestion window ($cwnd$).

Sulla base delle proprie informazioni locali, l'algoritmo modula la congestion window nel corso del tempo.

La modulazione provvede semplicemente ad incrementare la $cwnd$ (ogni versione usa tecniche differenti) finchè non si verifica o un evento di timeout, o un evento di loss.

Nel caso di un timeout si porta la $cwnd$ al valore minimo, in quanto si assume perso tutto ciò che è stato trasmesso (anche il BBR manterrà questo comportamento). Nel caso di un packet loss la $cwnd$ subisce un decremento moltiplicativo di un fattore β .

Quest'approccio non produce un controllo accurato del sending rate. Non è infatti possibile rilevare il momento in cui quest'ultimo eccede il bottleneck rate.

Il risultato conseguente sono condizioni operative altalenanti, in cui per la maggior parte si espandono i ritardi, e caricano i buffer intermedi.

Per essere più specifici, a seconda delle dimensioni dei buffer intermedi (in particolare del bottleneck buffer), potremmo riscontrare i seguenti problemi:

- *Shallow buffers* : in virtù della loro piccola dimensione, saranno molto frequenti i fenomeni di perdita. Una strategia AIMD porta a decrementi multipli continui della congestion window, quindi un delivery rate di piccola entità;
- *Deep buffers* : l'algoritmo reagirà nel momento in cui il bottleneck buffer sarà completamente saturo. L'RTT avrà raggiunto in questa condizione il suo picco massimo: valore che cresce all'aumentare della $BtlneckBufSize$. Quindi nel corso del tempo, è favorita la progressione dei ritardi, e delle perdite.

In aggiunta le perdite portano a dedicare parte della banda alle ritrasmissioni.

A questo punto è comprensibile la necessità di ricercare una strategia alternativa per limitare il sending rate, che non sia proibitiva in termini di loss rate, e che porti il delivery rate in prossimità del bottleneck rate.

Il controllo di congestione BBR

Il comportamento di BBR è dettato da un modello del path di rete, su cui il flusso di pacchetti viaggia. Tale modello traccia l'andamento nel corso del tempo del *BtlBw* e dell'*RTprop*.

Guidato dal proprio modello il BBR setta i propri parametri di controllo :

1. *pacing_rate*: è il parametro principale dell'algoritmo, e regola il sending rate del mittente nel corso del tempo;
2. *send_quantum*: definisce la massima unità di aggregazione per i pacchetti da trasmettere. Utile per diminuire l'overhead della CPU del destinatario.
3. *cwnd*: non ricopre un ruolo principale, ma continua ad essere il limite superiore per la quantità di dati inflight possibile del sender.

L'apposita configurazione di tali parametri permetterà al BBR di operare nel punto operativo migliore (target operating point). Ovvero quello in cui sono soddisfatte le condizioni di :

- *Rate balance* : il sending rate del mittente ed il bottleneck rate sono perfettamente bilanciati;
- *Full pipe* : la quantità di dati inflight è uguale al *BDP* corrente. La pipe opera al 100%, nulla più, nulla meno;

Nota: focus sull'installazione del BBR

Tutto ciò che successivamente diremo dell'algoritmo BBR, riguarda solo il lato mittente TCP. Non è necessario alcun cambiamento nel destinatario TCP per utilizzare quest'algoritmo di controllo di congestione.

E' importante però mettere l'accento sul fatto che ciò valga per il TCP che già fa uso di un sistema di ARQ (automatic repeat request).

Fondamentalmente il BBR potrebbe essere utilizzato anche da un altro protocollo di livello trasporto, l'importante è che disponga di un sistema di ARQ, in quanto gli ack sono alla base del corretto funzionamento dell'algoritmo.

4.1 Il network path model

In questa sezione approfondiamo come vengono misurati i parametri caratteristici del network path model.

4.1.1 Bottleneck bandwidth

Nella regione di bandwidth limited illustrata nel grafico 3.3, il delivery rate è praticamente pari al bottleneck rate, essendo questo il suo limite superiore. Quindi la stima del *BtlBw* può essere raggiunta attraverso il campionamento del delivery rate.

Il tempo con la quale sono generati i campioni utili alla stima è temporizzato dagli ACK ricevuti dal mittente. Approfondiamo adesso il processo di generazione dei campioni (per una trattazione specifica consultare [YCh17b]).

4.1.1.1 L'ack rate

Un ottimo indice possibile per il delivery rate è l'ack rate, ovvero il tasso di dati riscontrati istantaneamente dal destinatario.

Per il suo calcolo è necessario mantenere delle informazioni di stato, per ogni pacchetto, e per la connessione:

1. *C.delivered*: quantità di dati riscontrati, fino all'istante *C.delivered_time*;
2. *C.delivered_time*: istante in cui è pervenuto l'ack più recente;
3. *P.delivered*: valore di *C.delivered* all'atto della trasmissione di *P*;
4. *P.delivered_time*: valore di *C.delivered_time* all'atto della trasmissione di *P*;

A partire da queste è possibile definire l'espressione dell'ack rate:

$$ack\ rate = (C.delivered - P.delivered) / (C.delivered_time - P.delivered_time)$$

E' fondamentale comprendere che l'ack rate è vicino al delivery rate, se la quantità di dati riscontrata nel Δt individuato corrisponde a ciò che effettivamente il mittente trasferisce a destinazione.

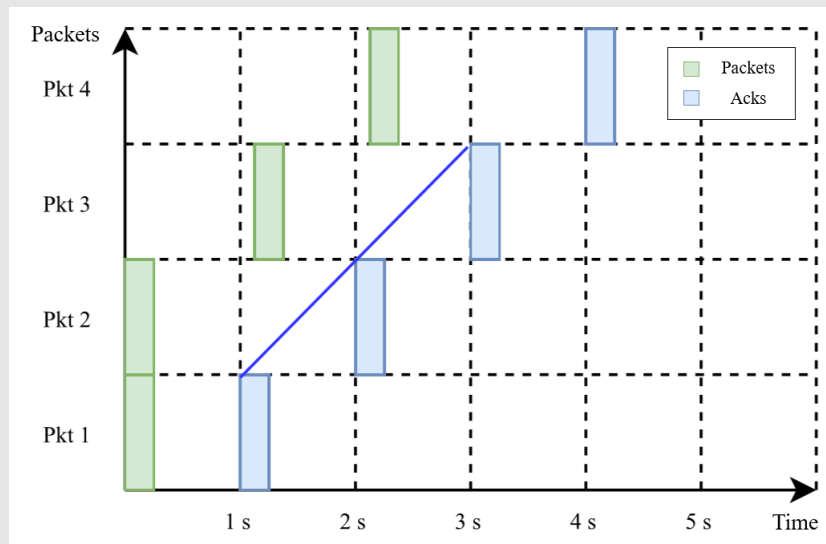
Fenomeni come la compressione, o soppressione degli ack potrebbero alterare la loro normale temporizzazione, portando a sovrastimare o sottostimare (rispettivamente) il delivery rate [Abr15].

La compressione è il principale nemico, in quanto una sovrastima per ragioni chiare avanti, è più pericolosa di una sottostima.

Approfondimento: significato matematico dell'ack rate

Mostriamo una possibile sequenza di trasmissione, da parte di un generico sender:

Figura 4.1: Ack rate determinato in corrispondenza del riscontro per Pkt 3



All'atto della ricezione del riscontro per il pacchetto 3, viene determinato l'ack rate come:

$$ack\ rate = (Pkt3.size + Pkt2.size) / (3s - 1s)$$

Dunque l'ack rate rappresenta la pendenza del segmento evidenziato. Ovvero la pendenza con la quale cresce il tasso di dati riscontrati.

In aggiunta si fa riflettere che il risultato (considerando pacchetti di dimensione fissa) è 1 pacchetto/secondo. Ed intuitivamente è ciò che ci saremo aspettati come delivery rate. (Ecco perchè l'ack rate è un buon indice in condizioni normali)

4.1.1.2 Il send rate

Per valutare la veridicità dell'ack rate viene introdotto il send rate.

Ancora una volta per la sua determinazione abbiamo bisogno di mantenere delle informazioni di stato (ulteriori alle precedenti), per ogni pacchetto:

1. `P.first_sent_time`: istante in cui è iniziata la trasmissione dell'ultimo pacchetto riscontrato dall'ack ricevuto al tempo `P.delivered_time` (a monte dell'invio di P)
2. `P.sent_time`: istante in cui parte la trasmissione del pacchetto

A partire da queste definiamo l'espressione del send rate:

$$send\ rate = (C.delivered - P.delivered) / (P.sent_time - P.first_sent_time)$$

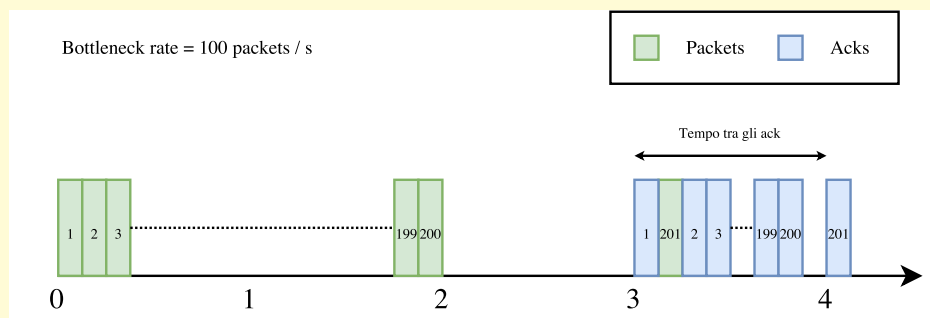
La differenza con l'ack rate sta nell'intervallo di tempo considerato. Il send time rappresenta il tempo che è trascorso tra la trasmissione dell'ultimo pacchetto riscontrato, ed il pacchetto attualmente riscontrato.

Questo permette che il Δt considerato includa il tempo di trasmissione dei dati riscontrati, in quanto il send time è in generale \geq dell'effettivo tempo di trasmissione. Così è possibile evitare il problema delle sovrastime dell'ack rate.

Esempio: un caso di compressione degli ack

Mostriamo una possibile sequenza di trasmissione, in cui si verifica una compressione degli ack ricevuti:

Figura 4.2: Fenomeno di compressione degli ack



Come si nota dalla figura gli ack per i pacchetti 2-200 vengono ritardati e compressi, ed arrivano in burst tra la trasmissione del pacchetto 201 ed il relativo ack.

Ciò porta l'ack rate ad essere 200 packets/s, ma si comprende bene che tra 3s e 4s il sender non ha trasferito 200 pacchetti a destinazione. Quindi in tale situazione l'ack rate non è un indice veritiero del delivery rate.

Il reale tempo che i 200 pacchetti impiegano per arrivare al destinatario è di 2s.

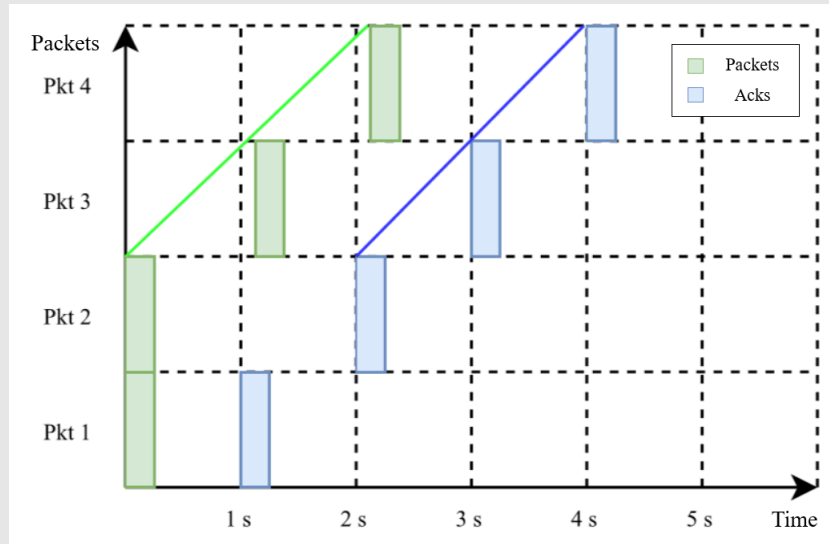
In questo caso il send rate è di 60 packets/s, perchè considera il tempo trascorso dalla trasmissione del pacchetto 1 a quella del pacchetto 201. Il secondo in più è dovuto all'inclusione del secondo tra 2 e 3 secondi, in cui il mittente è inattivo.

Confrontando quindi il send rate e l'ack rate si comprende che il destinatario non può aver ricevuto 200 pacchetti tra 3s e 4s.

Approfondimento: significato matematico del send rate

Consideriamo la stessa sequenza di trasmissione precedente da parte di un generico sender:

Figura 4.3: Send rate determinato in corrispondenza del riscontro per Pkt 4



All'atto della ricezione del riscontro per il pacchetto 4 viene determinato il send rate come :

$$send\ rate = (Pkt3.size + Pkt4.size) / (2s - 0s)$$

Dunque il send rate rappresenta la pendenza del segmento evidenziato in verde. Ovvero la pendenza con la quale cresce il tasso di dati trasmessi.

4.1.1.3 Stima del *BtlBw*

Ad ogni ack sarà generato un campione relativo al delivery rate come segue:

Codice 4.1: Delivery rate evaluation

```
time_elapsed = max(send_time , ack_time)
delivery_rate = data_acked / time_elapsed
```

Non ci resta altro che scoprire come da tali campioni venga ottenuta una stima del bottleneck rate.

Per poter far questo BBR usa un *BtlBwMaxFilter*, un windowed filter di ampiezza W_B :

$$\widehat{BtlBw}_T = \max(deliveryRate_t) \forall t \in [T - W_B, T]$$

L'ampiezza W_B è misurata in multipli di $RTprop$, e deve essere:

- Abbastanza estesa da contenere un ProbeBW cycle gain (sezione 4.3.3). Così che essa contenga dei sample raccolti nella regione di bandwidth limited, che è effettivamente l'unica regione in cui la stima potrà essere corretta;
- Abbastanza estesa da combattere l'eventuale presenza di rumore, introdotta da eventuali code lungo il percorso;
- Abbastanza corta da poter permettere una pronta reazione ad una diminuzione del bottleneck rate;

Un ottimo valore è costituito da 10 $RTprops$ [Car+17b, p. 61].

4.1.2 Round Trip propagation delay

Per la stima dell' $RTprop$ siamo aiutati dal fatto che il protocollo TCP già effettua un campionamento dell' RTT , per la gestione dei timeout di ritrasmissione.

Anche per determinare $RTprop$ BBR usa un windowed filter, $RTpropMinFilter$ di ampiezza W_R :

$$\widehat{RTprop}_T = \min(RTT_t) \forall t \in [T - W_R, T]$$

L'ampiezza W_R è misurata in secondi, e deve essere :

- Estesa almeno quanto un ProbeRTTInterval (sezione 4.3.4). Così che essa contenga dei sample raccolti nella regione di application limited, che è effettivamente l'unica regione in cui la stima potrà essere corretta;
- Abbastanza estesa da evitare una consistente diminuzione del delivery rate raramente (sezione 4.3.4);
- Abbastanza estesa da poter sfruttare gli istanti in cui l'applicazione limita il delivery rate del mittente, per poter aggiornare la stima dell' $RTprop$ senza diminuire fortemente il delivery rate;
- Abbastanza corta da poter permettere una pronta reazione ad una variazione dei ritardi;

4.2 I parametri di controllo

Compreso in primo luogo come il BBR mantenga aggiornato il modello del network path. Vediamo ora di approfondire i parametri di controllo, usati da quest'ultimo per portarsi nel target operating point.

4.2.1 Pacing rate

Il pacing rate del BBR costituisce un parametro importantissimo (e primario), con cui vengono incontrate nel tempo le condizioni di *Rate balance*, e di *Full pipe*.

Effettuare il pacing dei pacchetti significa equi-distanziare le loro trasmissioni. Quindi a monte dell'invio di un dato pacchetto il successivo sarà trasmesso al tempo:

$$nextSendTime = now + packet.size / pacing\ rate$$

Ciò porta ad un adattamento del tasso di trasmissione del mittente al pacing rate, nonostante la capacità del proprio link d'ingresso possa essere superiore. La spaziatrice definisce dunque le pause adatte per sincronizzare il sending rate come voluto.

Il set di tale parametro (così come gli altri) avviene ad ogni ack, e può essere riassunto con il seguente pseudocodice:

Codice 4.2: BBRSetPacingRate

```

BBRSetPacingRate () :
2   BBRSetPacingRateWithGain (BBR.pacing_gain )

4   BBRSetPacingRateWithGain ( pacing_gain ) :
    rate = pacing_gain * BBR.BtlBw
6   if (BBR.filled_pipe || rate > BBR.pacing_rate )
    BBR.pacing_rate = rate

```

Da cui come ci si aspettava il pacing rate è settato proporzionalmente al *BtlBw*, attraverso il pacing gain. La condizione dell'if non ne permette il decremento in condizioni di pipe underfilled.

Un valore di pacing gain pari 1 può permettere di operare nel target operating point, se non sono state formate code in precedenza. Un valore di pacing gain > 1 è utile per ottenere sample utili alla stima del *BtlBw*, in quanto viene saturato il delivery rate. Ed infine un valore di pacing gain < 1 può essere utilizzato per svuotare le code generate, ed uscire da una situazione di pipe overfill.

4.2.2 Send quantum

Rappresenta la massima unità di aggregazione per i pacchetti inviati dal mittente. Permette quindi d'inviare gruppi di pacchetti in un unico burst, diminuendo l'overhead della CPU del destinatario.

Settato ad ogni ack, il suo valore dipende fortemente dal pacing rate. Per valori alti di quest'ultimo è possibile usare alti valori di send quantum, mentre per valori bassi è meglio diminuirlo, in quanto creerebbe un aumento dei ritardi e delle code lungo il percorso.

4.2.3 Congestion window

La *cwnd* impone un limite alla quantità di dati in flight del mittente. Nonostante il parametro principale del BBR per il controllo del sending rate sia il pacing rate, la *cwnd* è molto utilizzata nelle situazioni in cui si vuole subito portare il limite in flight in un particolare punto. Verrebbe speso un maggior tempo se si usufruisse del pacing rate.

Per la sua gestione vengono definiti un upper ed un lower bound:

- *target cwnd*: limite superiore della *cwnd*, ottenuto in funzione del \widehat{BDP} (*BDP* stimato) come:

$$target\ cwnd = cwnd\ gain * \widehat{BDP}$$

- *BBRMinPipeCwnd*: limite inferiore della *cwnd*, tipicamente tale da consentire l'invio di 2/4 pacchetti al mittente.

Tali limiti sono da considerare strettamente rigidi in condizioni normali (non occorrono nè eventi di loss, nè di timeout).

Discutiamo adesso della modulazione del parametro in tali condizioni.

Se la pipe fosse piena, è molto probabile che si sia raggiunto il target, dunque si provvede ad incrementare la *cwnd* solo se essa non supererà il target. In caso contrario resterà fissata su quest'ultimo.

Se invece fossimo al di sotto del target, la *cwnd* subirebbe un semplice incremento pari alla quantità di dati consegnati a destinazione.

In ogni caso viene sempre garantito il rispetto del lower bound.

Tuttavia gli eventi di packet loss e di timeout portano ad una modulazione diversa della *cwnd*.

Ad un evento di timeout, il BBR porta la *cwnd* cautamente ad una SMSS, in quanto viene assunto perso tutto ciò che è stato inviato. Da qui poi si ripartirà come descritto in precedenza.

In occorrenza invece di un evento di packet loss, la *cwnd* viene subito (al primo round) adeguata al delivery rate attuale come segue:

Codice 4.3: BBR enter in loss recovery

```

2 BBR.prior_cwnd = BBRSaveCwnd()
  cwnd = packets_in_flight + max(packets_delivered, 1)
4 BBR.packet_conservation = true

```

Adeguandosi poi, nelle fasi successive, sempre al delivery rate attuale:

Codice 4.4: BBRModulateCwndForRecovery

```

1 BBRModulateCwndForRecovery() :
3   if (packets_lost > 0)
        cwnd = max(cwnd - packets_lost, 1)
5   if (BBR.packet_conservation)
        cwnd = max(cwnd, packets_in_flight + packets_delivered)

```

Dopo un round trip time in Loss Recovery, si ripristina la cwnd al suo precedente valore (salvato con BBRSaveCwnd):

Codice 4.5: BBR exit from loss recovery

```

1 BBR.packet_conservation = false
  BBRRestoreCwnd()
3
  BBRRestoreCwnd() :
5   cwnd = max(cwnd, BBR.prior_cwnd)

```

A questo punto riportiamo il codice completo di gestione della congestion window:

Codice 4.6: BBRSetCwnd

```

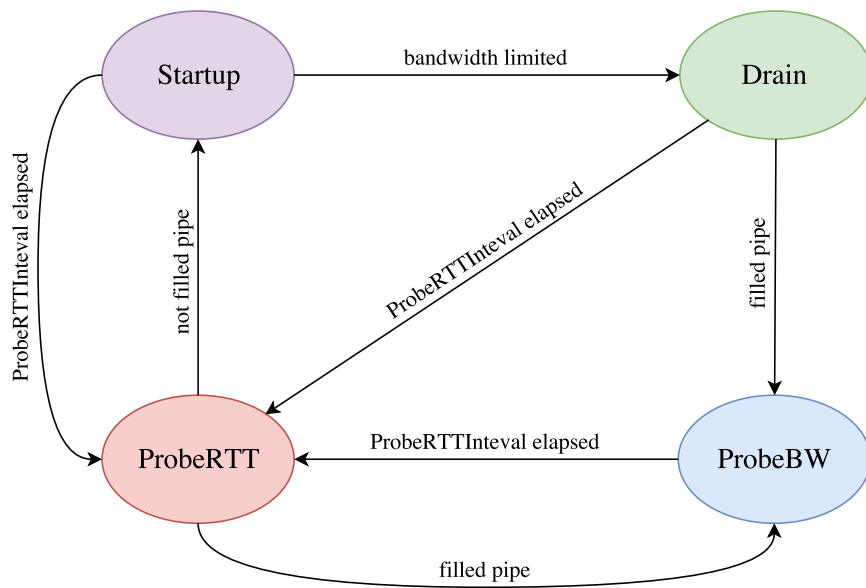
  BBRSetCwnd() :
2   BBRUpdateTargetCwnd()
  BBRModulateCwndForRecovery()
4   if (not BBR.packet_conservation) {
        if (BBR.filled_pipe)
6           cwnd = min(cwnd + packets_delivered, BBR.target_cwnd)
        else if (cwnd < BBR.target_cwnd || BBR.delivered <
            InitialCwnd)
8           cwnd = cwnd + packets_delivered
        cwnd = max(cwnd, BBRMinPipeCwnd)
10    }
  BBRModulateCwndForProbeRTT()

```

4.3 Visione complessiva dell'algoritmo

In questa sezione fonderemo tutti gli aspetti visti in precedenza. Comprendendo affondo come lavora un flusso BBR nel corso della sua vita.

Il miglior punto di partenza è visualizzare gli stati in cui il flusso può trovarsi con uno state chart diagram:

Figura 4.4: BBR state diagram

Ricordiamo che gli obiettivi della state machine sono:

- Raggiungere il punto di massimo per il delivery rate, ed il punto di minimo per l'RTT;
- Mantenere continuamente aggiornato il network path model, stimando sequenzialmente nel corso del tempo, il bottleneck rate ed il round trip propagation delay;
- Garantire una equa suddivisione della capacità di rete tra flussi BBR concorrenti;

Passiamo ora all'analisi degli stati e delle transizioni.

4.3.1 Startup

Durante la fase di Startup, così com'era per il classico Slow Start, è d'interesse scovare subito il limite massimo del delivery rate.

Ora considerando gli standard attuali è ragionevole avere a che fare con 9 ordini di grandezza (1 Gbps), 10 ordini o 11 ordini di grandezza (10/100 Gbps).

Pertanto è necessario che progressivamente ad ogni round, il sending rate perlomeno raddoppi. Ciò viene ottenuto configurando il pacing gain, ed il cwnd gain ad un valore pari al BBRHighGain ($2/\ln(2)$). Ovvero la quantità di guadagno minima che garantisce un raddoppio del tasso del mittente:

Codice 4.7: BBREnterStartup

```

BBREnterStartup() :
2   BBR.state = Startup
   BBR.pacing_gain = BBRHighGain
4   BBR.cwnd_gain = BBRHighGain
  
```

La fase di Startup termina nel momento in cui si riempie completamente la pipe. Ciò che tipicamente accade, è che non si riesce ad arrivare precisamente nel punto operativo, ma un pò aldilà, generando così delle code di una quantità pari a $(\text{cwnd_gain} - 1) * BDP$.

Per individuare poi in che momento si è saturato il delivery rate, viene usato il network path model. Ciò che si fa è attendere che per 3 round di fila il bottleneck rate stimato non cresca più del 25%. Quindi viene ricercato fondamentalmente un plateau per quest'ultimo.

Motiviamo poi perchè siano necessari 3 round. Il primo round serve al destinatario per effettuare l'auto-tuning della finestra di ricezione, sulla base dell'attuale capacità di rete. Il secondo round è usato dal mittente per riempire la finestra di ricezione. Il terzo è quello in cui arriveranno i riscontri per i pacchetti inviati, e cioè i sample che effettivamente sono indicativi del delivery rate.

4.3.2 Drain

Dopo una fase di Startup entriamo in Drain. Lo scopo di tale fase è svuotare le code formate in precedenza. Ciò viene ottenuto in un round invertendo il valore del pacing gain (così da lavorare con un sending rate $<$ bottleneck rate):

Codice 4.8: BBREnterDrain

```

1 BBREnterDrain() :
  BBR.state = Drain
3 BBR.pacing_gain = 1/BBRHighGain
  BBR.cwnd_gain = BBRHighGain

```

La fase di drain termina nel momento in cui la quantità di dati inflight arriva al valore attuale stimato del BDP . Quindi le code saranno tutte svuotate.

A questo punto si passa in ProbeBW.

4.3.3 ProbeBW

Il BBR arriva nello stato ProbeBW senza la presenza di code sul path di comunicazione.

E' in questo stato in cui il BBR spende la maggior parte del suo tempo.

In ProbeBW il BBR stima il bottleneck rate, operando in condizioni di massimo delivery rate e ritardo di propagazione minimo.

Ciò viene ottenuto attraverso una tecnica chiamata gain cycling. Un gain cycle è tipicamente costituito da 8 fasi, in ognuna di queste fasi il pacing gain viene impostato ad un valore utile a raggiungere le condizioni sopra espresse.

I valori di pacing gain usati per ogni fase sono i seguenti: $[5/4, 3/4, 1, 1, 1, 1, 1, 1]$.

Partiamo dalla fase 0. In questa fase il pacing gain porta il mittente ad operare con un pacing rate $> BtlBw$. Questo serve ad ottenere campioni utili alla stima di quest'ultimo, in quanto è solo nella regione di bandwidth limited (figura 3.3) ove il delivery rate \simeq bottleneck rate.

La fase 0 dura almeno un $RTprop$, e termina nel momento in cui la quantità di dati inflight è pari a $5/4 * BDP$, o si verificano degli eventi di loss (non è possibile raggiungere la condizione inflight $\geq 5/4 * BDP$).

Dopo la fase 0, l'obiettivo della fase 1 è drenare le code che sono state formate, impostando un valore di pacing rate di $3/4$. Dunque in un tempo pari a $RTprop$, o prima (se le code si svuotano rapidamente) il flusso BBR torna nuovamente nel target point.

Dopodichè dalla fase 2 in poi, quest'ultimo opera nelle condizioni ideali con un pacing gain pari a 1.

Proponiamo di seguito uno scenario reale di funzionamento del cycle gaining, in cui il bottleneck rate è di 10Mbps:

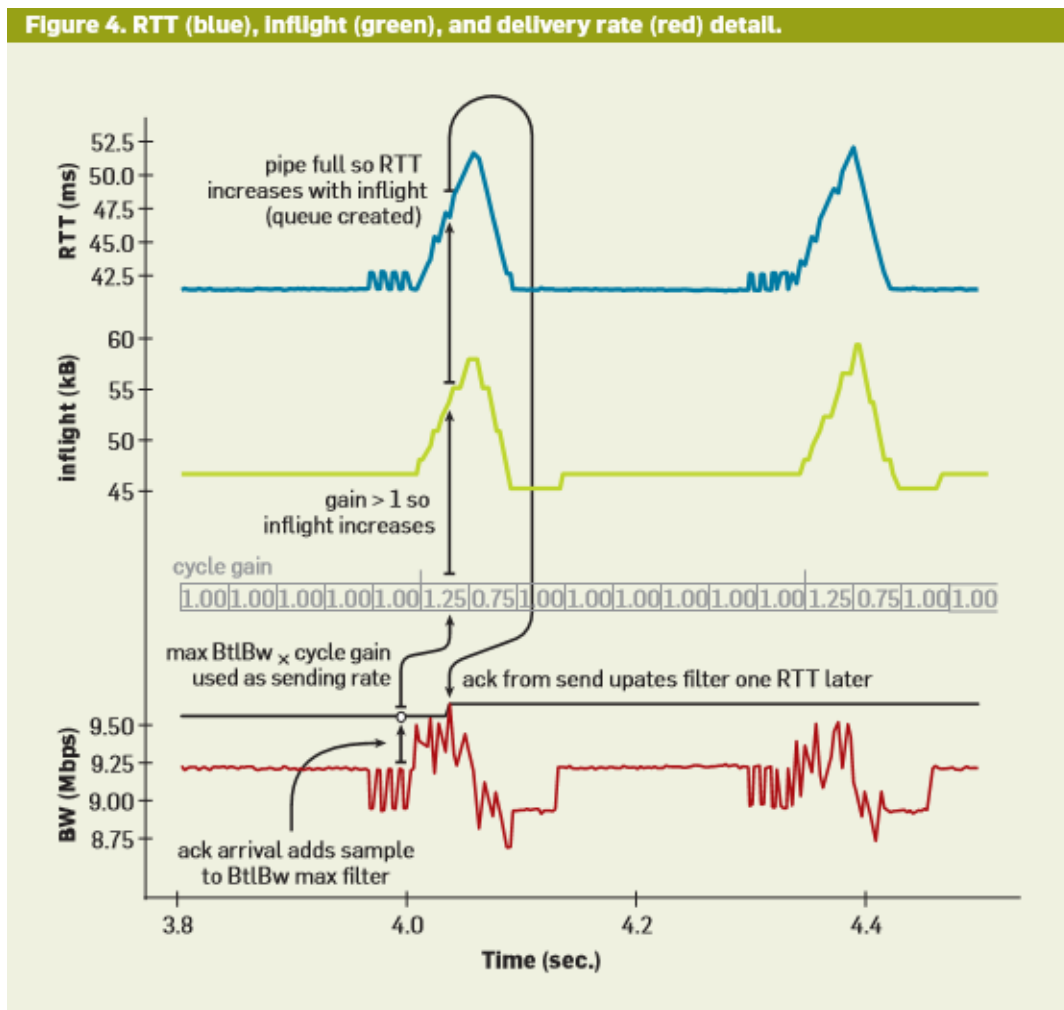
Figura 4.5: BBR probe bw gain cycling

Figura tratta da [Car+17b, p. 62]

Aggiungiamo solo che in questo stato la *cwnd* non ha un ruolo centrale, ed il suo target viene impostato a $2 \cdot BDP$. Questo per non disturbare la modulazione del pacing rate, e per avere comunque un limite di sicurezza per l'inflight amount.

4.3.4 ProbeRTT

Il BBR usa tale stato per poter ottenere sample dell'RTT utili a stimare correttamente l'*RT_{prop}*.

Ogni ProbeRTTInterval (di una durata tipica di 10 s) il BBR transita in ProbeRTT se la stima non è stata aggiornata. In tale stato ci si serve della *cwnd* per dare un colpo netto alla quantità di dati inflight, impostandola al suo limite minimo (BBRMinCwndPipe). In questo modo il flusso opera nella regione di application limited, ove si otterranno sample dell'RTT $\simeq RT_{prop}$.

La *cwnd* è mantenuta a *BBRMinCwndPipe* per un tempo pari al *ProbeRTTDuration*, che è circa 200 ms. Questa scelta permette di spendere solo il 2% del tempo in *ProbeRTT*, dove il delivery rate è limitato fortemente, a discapito del 98% speso giustamente in *ProbeBW*.

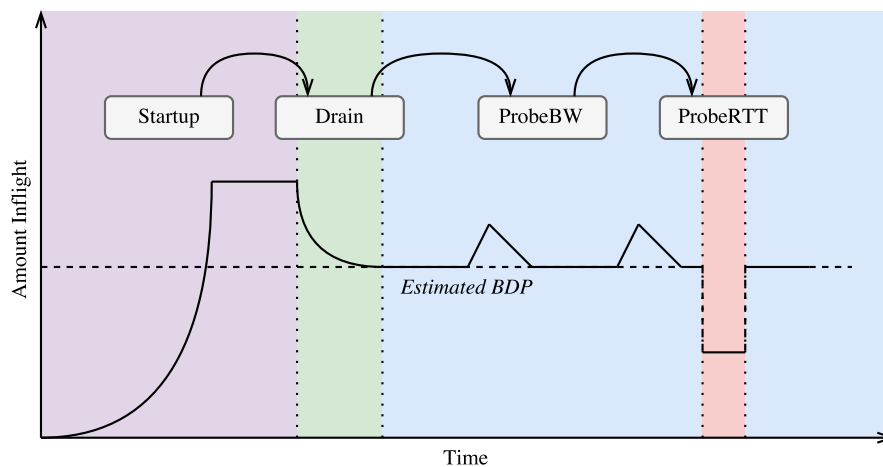
Allo scadere di un *ProbeRTTDuration* viene ripristinata la *cwnd*, e si transita o in *ProbeBW*, o in *Startup*, a seconda dello stato attuale della rete.

4.3.5 Evoluzione di un flusso BBR

In conclusione possiamo dire che l'evoluzione di un flusso BBR, passa per un:

1. *Warm up*: si passa in stati non stabili come *Startup* e *Drain* per scovare velocemente il limite della banda di trasmissione;
2. *Stati stabili*: quali *ProbeBW* e *ProbeRTT*, in cui il BBR aggiorna il proprio network path, ed imposta i propri parametri per raggiungere il target operating point;

Figura 4.6: BBR state evolution



Casi reali d'impiego del BBR

5.1 Cubic vs BBR

Sia Cubic (App. A) che BBR, si servono di uno stato iniziale, quali lo Slow Start e lo Startup (rispettivamente), per scovare rapidamente il limite superiore del delivery rate.

Tuttavia, hanno un intento differente, essendo due approcci al controllo di congestione diversi.

Il Cubic cerca il limite superiore per fissare la soglia di passaggio dallo Slow Start, al Congestion Avoidance. In cui semplicemente, continuerà ad aumentare il sending rate, ma in modo più lento. Questo finchè non raggiungerà il punto in cui scadrà un timeout, o il bottleneck buffer sarà saturo. Condizione in cui si verificheranno delle perdite.

Il BBR invece, raggiunto il limite del delivery rate, con la consapevolezza di esser andato un pò oltre (il target point) e di aver generato delle code, passa in Drain per scaricarle. Quindi arriva in ProbeBW con le condizioni giuste e necessarie ad operare regolarmente al massimo tasso, ed ai minimi ritardi.

Possiamo infatti sottolineare la loro evidente differenza nella gestione dell'RTT, con il seguente scenario:

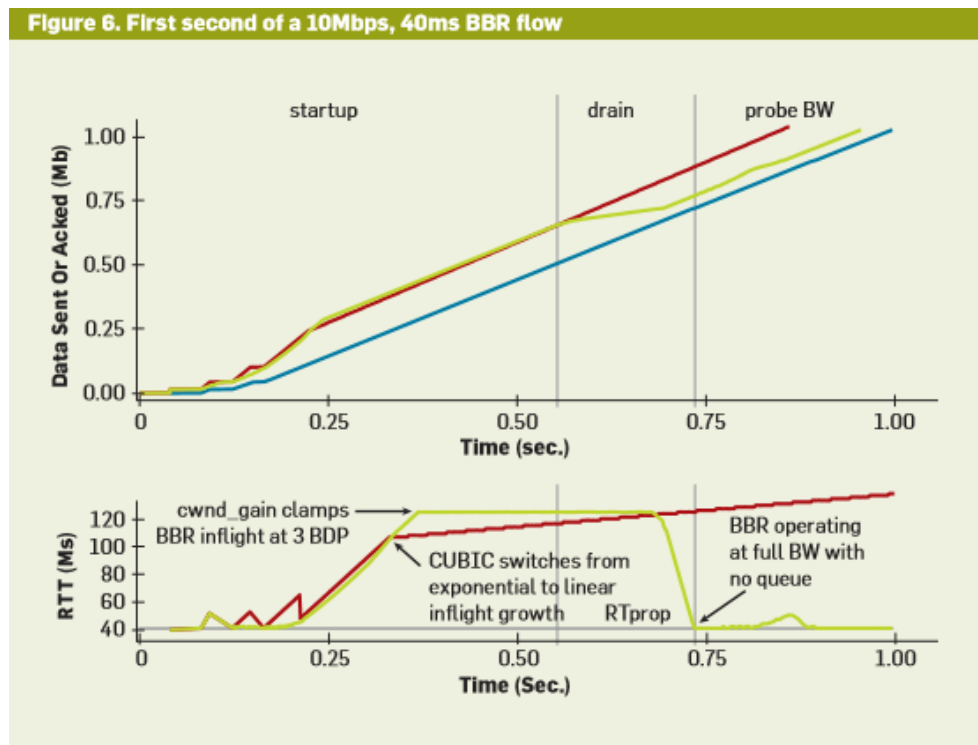
Figura 5.1: First second of a 10Mbps, 40ms BBR flow

Figura tratta da [Car+17b, p. 63]

L'ambito operativo di Cubic è sicuramente quello peggiore, in quanto avremo dei livelli alternati di RTT, tra il 70% ed il 100% (con 0% intendiamo la condizione in cui $RTT \approx RTprop$), ed un delivery rate si al massimo, ma dimezzato ogni volta in corrispondenza di una perdita (o abbattuto ad un timeout):

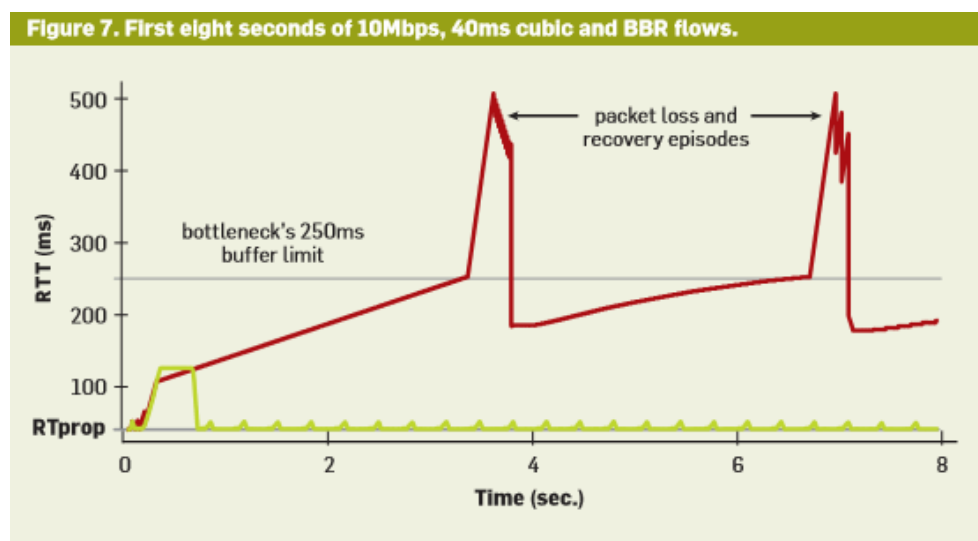
Figura 5.2: First eight seconds of a 10Mbps, 40ms cubic and BBR flow

Figura tratta da [Car+17b, p. 64]

Bisogna però prendere atto di una cosa, ovvero che Cubic a differenza del BBR non dispone di un network path model. Quindi non è possibile contenere l'RTT se non lo si traccia nel tempo, e se non è quest'ultimo a determinare il valore dei parametri di controllo della soluzione.

5.2 Fairness tra flussi BBR

Tra flussi BBR concorrenti, il livello di fairness raggiunto, testimoniato anche da appositi esperimenti, è molto soddisfacente.

Ecco infatti un caso d'esempio di 5 flussi concorrenti, che condividono un bottleneck link da 100 Mbps, ed un RTT da 10 ms:

Figura 5.3: Throughput of five BBR flows sharing a bottleneck

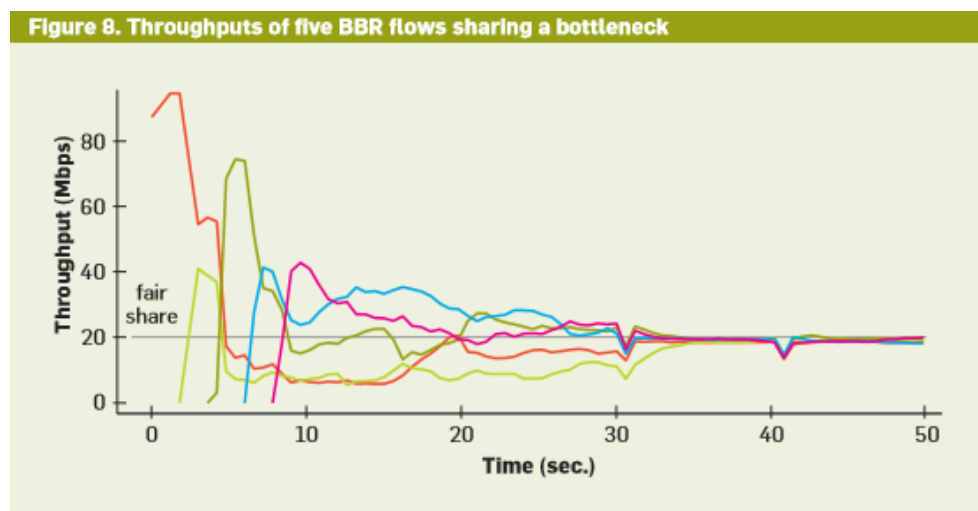


Figura tratta da [Car+17b, p. 64]

Nella figura viene presentato uno scenario abbastanza generale, in cui i flussi BBR non nascono tutti nello stesso istante.

Ciò che ci aspettiamo accada è semplice. Nello stato di ProbeBW i flussi al di sopra del fair share scopriranno un calo del bottleneck rate, quindi si accorgeranno della perdita di banda, riconfigurando poi il loro pacing rate.

Andando poi avanti, quando (di nuovo) un flusso al di sopra del fair share - che dunque sta generando un maggior traffico - entrerà in ProbeRTT, svuoterà considerevolmente il bottleneck buffer. Ciò porterà ad un aggiornamento dell'RTprop degli altri flussi concorrenti.

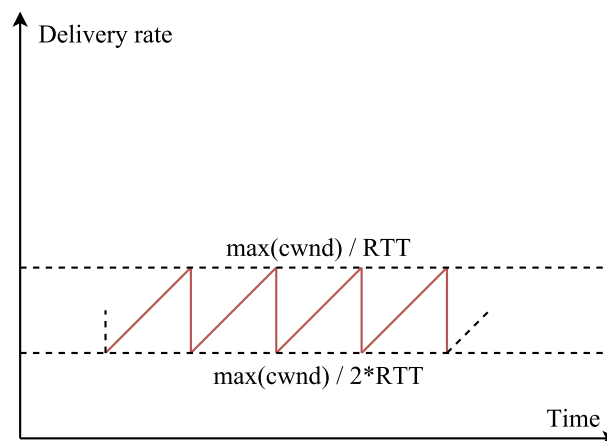
Alla lunga, il tutto porta a sincronizzare l'entrata nello stato di ProbeRTT di tutti i flussi. Così tutti potranno stimare il reale valore dell'RTprop, così come del bottleneck rate.

5.3 Fairness tra flussi loss-based e BBR

In questo caso, è abbastanza (purtroppo) prevedibile che allo stato attuale delle cose, sia difficile raggiungere una suddivisione equa della banda.

Fondamentalmente, un flusso loss-based presenta un andamento del throughput pressochè a dente di sega:

Figura 5.4: TCP loss-based evoluzione del throughput



Allora ci aspetteremo che un flusso BBR abbia maggiore banda nelle fasi di low-throughput del flusso loss, mentre questi diminuisca quando quest'ultimo arriva al suo massimo (le sue reali possibilità). Tuttavia, un flusso loss, ad esempio come Reno, raggiunge il massimo in corrispondenza di un evento di perdita, per cui non manterrà a lungo tale livello.

Per la consultazione di appositi esperimenti si veda [Car+17a].

5.4 Rete Google B4

B4 è una rete proprietaria ad estensione geografica, dalle alte performance, di Google.

Nel 2015, quando BBR non era più solo una teoria, ma una idea concreta, si diede inizio ad una transizione verso quest'ultimo. Così nel 2016, Cubic è stato completamente sostituito dal BBR, che non essendo un approccio loss-based, ha portato un notevole incremento delle prestazioni.

Il primo artefatto d'analisi mette a rapporto il throughput raggiunto con BBR rispetto quello raggiunto con Cubic:

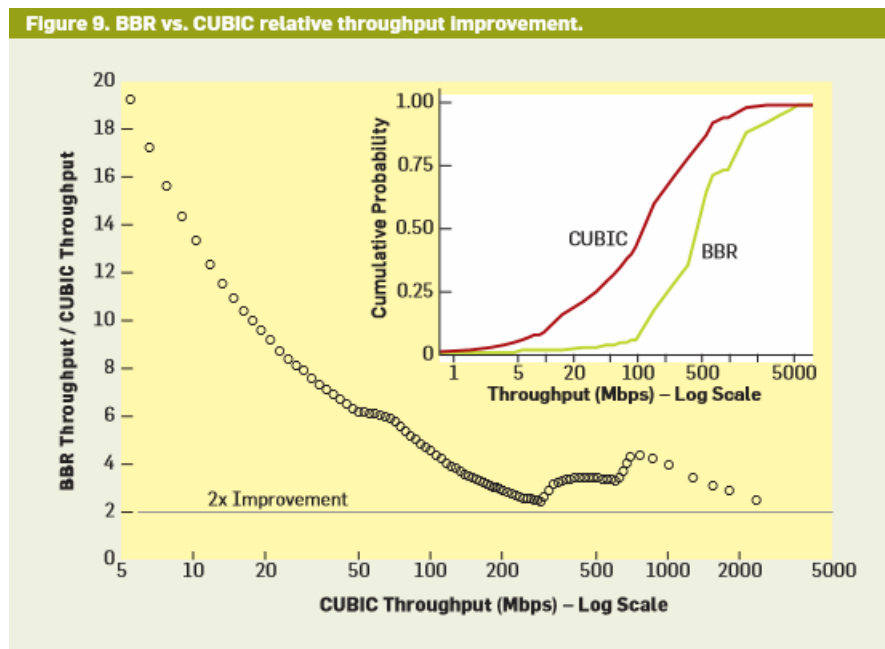
Figura 5.5: BBR vs CUBIC relative throughput improvement

Figura tratta da [Car+17b, p. 64]

Tali dati, raccolti da varie sonde disposte in vari punti della rete B4, mettono in evidenza un incremento prestazionale che va dal 2x al 25x. In realtà, l'incremento del BBR è mascherato dal limite imposto ai buffer di ricezione TCP (8 MB).

Alzando tale limite, è stato ottenuto su un path tra Europa e USA, un throughput di 2 Gbps, contro i 15 Mbps di Cubic. Parliamo di un incremento: 133x.

Nel grafico vengono anche messe a confronto le distribuzioni di probabilità cumulative. Da esso possiamo notare (qualitativamente) che:

Tabella 5.1: Throughput CDF

	BBR	Cubic
$P(\text{thput} < 1)$	$\simeq 0$	$\simeq 0$
$P(\text{thput} < 20)$	$\simeq 0$	0.19
$P(\text{thput} < 100)$	0.06	0.5
$P(\text{thput} < 500)$	0.62	0.87
$P(\text{thput} < 5000)$	1	1

Il secondo artefatto mette a confronto il throughput utile sostenibile in funzione del loss-rate:

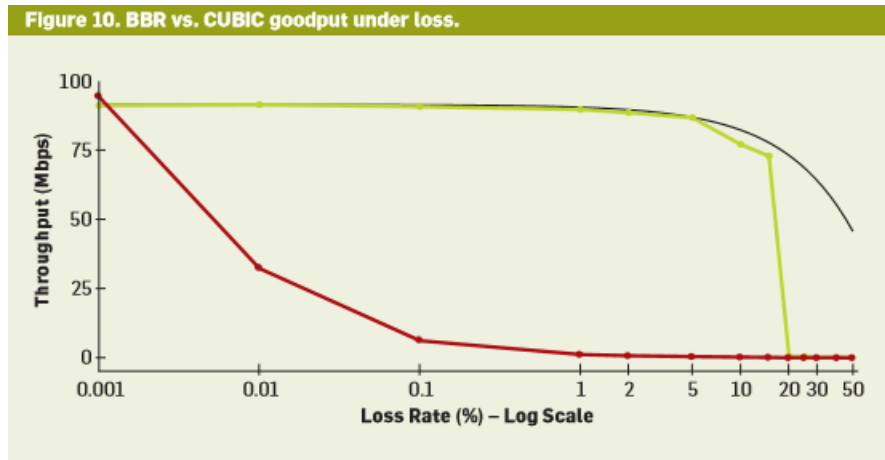
Figura 5.6: BBR vs CUBIC goodput under loss

Figura tratta da [Car+17b, p. 65]

Ricordando che un loss-rate ragionevole è circa dell'1%, e notando che il BBR perde performance tra 5% ed 20%, questo è un ottimo miglioramento.

Appendice A

Nel corso dell'elaborato si è fatto riferimento alla versione Cubic del TCP, come metro di giudizio per le prestazioni di BBR. In quest'appendice si richiamano alcuni aspetti del Cubic, utili a chiarire alcuni dei suoi principali comportamenti.

Tale appendice non è da considerarsi come fonte principale per lo studio di Cubic. Per una trattazione approfondita si suggerisce la lettura di [RSc17], e [HRX08].

6.1 Motivazione

La motivazione principale che ha portato alla nascita di Cubic, riguarda un problema presentato dalle versioni Reno, e New Reno, su path di comunicazione caratterizzati da livelli elevati di RTT.

Tali algoritmi, durante la fase di *Congestion Avoidance*, eseguono un incremento della `cwnd` di una SMSS per ogni RTT. Questo causa un lento incremento della finestra, per path caratterizzati da alti valori del *BDP*.

Per esempio, con un bottleneck rate di 10Gbps e 100ms di RTT, il *BDP* è di 100.000 pacchetti, se questi hanno una dimensione di 1250 Byte. Questo significa che per incrementare la `cwnd`, da *BDP/2* (il punto di partenza è lo *Slow start*, non il *Congestion Avoidance*) fino a *BDP*, è necessario un tempo pari a: $50.000 * RTT = 5.000 \text{ s}$ (1.4 ore).

Tutto ciò non porta altro che ad una forte sottovalutazione delle capacità reali del path di rete.

Sarebbe dunque utile, essere più aggressivi nella fase d'incremento quando si è lontani dal punto di saturazione, e meno nelle sue vicinanze.

Il Bic-TCP fu una, tra le tanti versioni, a riuscire in tale intento.

6.2 Bic-TCP

Il Bic-TCP è un algoritmo di controllo di congestione ottimizzato per i path con alti valori del *BDP*.

Per raggiungere tale risultato, durante la fase di *Congestion Avoidance*, viene utilizzato un algoritmo di binary search.

Per utilizzarlo è innanzitutto necessario nel corso del tempo, documentare i limiti della finestra di congestione:

- W_{max} : ultimo valore della $cwnd$ in cui si è generato un evento di packet loss, registrato prima d'entrare in *Fast Recovery*.
- W_{min} : ultimo valore della $cwnd$ in cui per un intero RTT, non si è verificato un evento di packet loss.

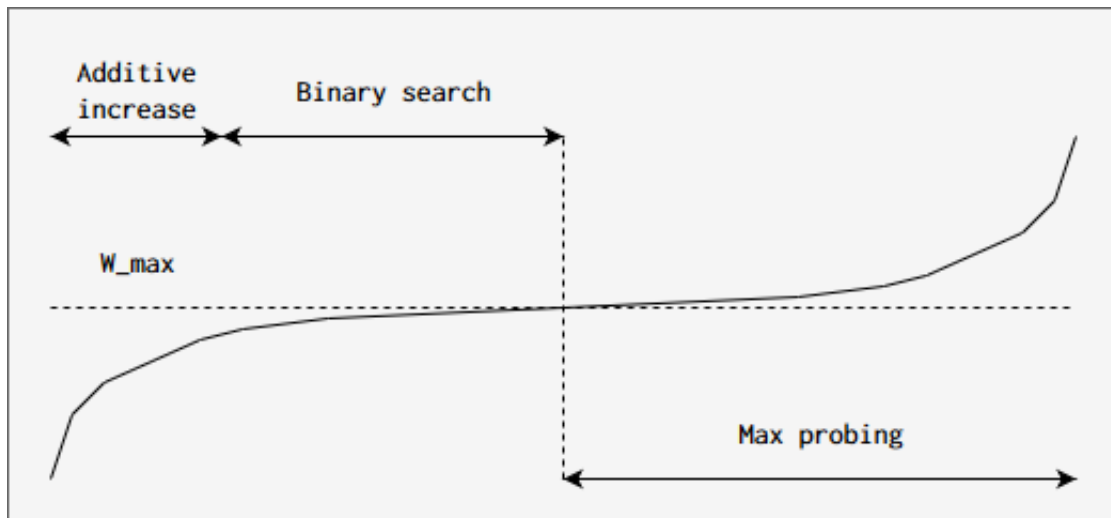
A questo punto in *Congestion Avoidance*, ad ogni ack ricevuto:

- Se $cwnd < W_{max}$, viene realizzata una binary search:
 1. viene determinato il mid-point tra W_{min} (pari al valore corrente di $cwnd$) e W_{max} .
 2. se non è possibile in un RTT incrementare la $cwnd$ fino al mid-point, viene eseguito un incremento addittivo di un fattore costante S_{max} (parametro di configurazione). In caso contrario l'incremento è dato dalla distanza rispetto il mid-point.
 3. se per un RTT non si presentano eventi di loss, viene aggiornato il valore di W_{min} , con l'attuale valore della $cwnd$. In caso contrario si riprenderà dal passo 1, all'uscita dalla fase di *Fast Recovery*.
 4. Finchè l'incremento non scende al di sotto S_{min} (parametro di configurazione, che segnala se si è o meno nelle vicinanze del punto di saturazione), si riparte dal passo 1.

Nel suo insieme il binary search, dà vita ad un incremento pressochè logaritmico, molto rapido e stabile in prossimità del punto di saturazione.

- Se $cwnd \geq W_{max}$, significa che il limite della finestra congestione ha subito una variazione. Per cui si procede con una fase d'incremento duale alla precedente, prima esponenziale e poi lineare (Max probing).

Mettendo insieme il tutto, otteniamo l'andamento della funzione di crescita del Bic (growth function):

Figura 6.1: Bic growth function

La forza del Bic, può essere desunta dal plateau che si forma nell'intorno di W_{max} . Quest'ultimo permette di operare con valori alti della $cwnd$ (prossimi al limite), per un tempo superiore agli approcci tradizionali.

Se dunque è comprensibile, che questo modo di modulare la finestra di congestione, aiuta quando il BDP in gioco è elevato, è altrettanto comprensibile che un tale approccio sia sicuramente poco fair nei confronti di Reno per BDP bassi.

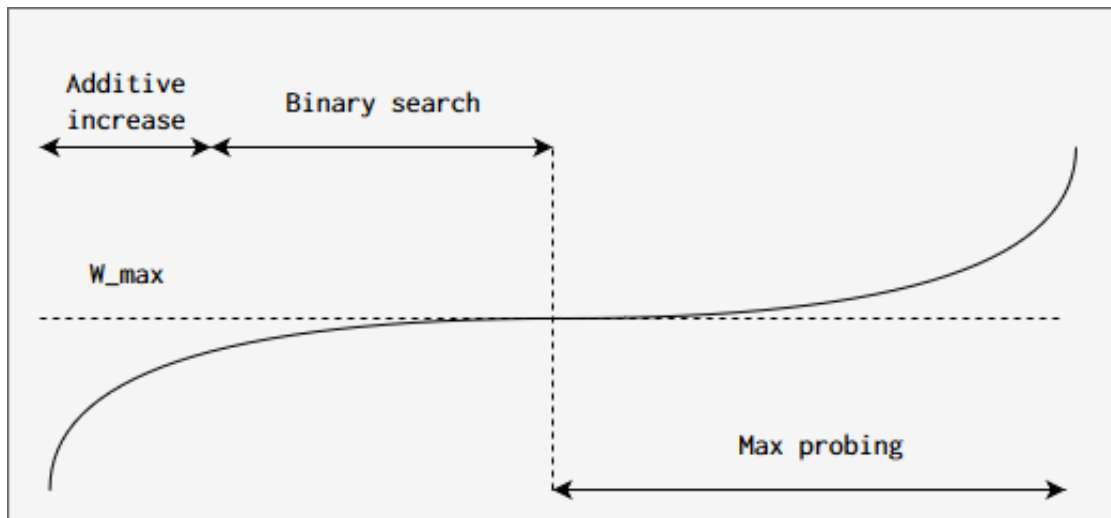
Non solo, è possibile che in alcune zone di funzionamento la Bic function produca un avanzamento della finestra di congestione più lento rispetto quello lineare. Ciò può accadere quando si è lontani dal punto di saturazione.

Tutto ciò congiuntamente al fatto che la realizzazione di un binary search algorithm, non è cosa semplice, ha portato allo sviluppo di Cubic.

6.3 Cubic

Cubic è un'evoluzione del Bic-TCP, di cui uno degli scopi è proprio quello di migliorare gli aspetti negativi di quest'ultimo.

Non è infatti basato sull'utilizzo del binary search algorithm, ma sulla Cubic function:

Figura 6.2: Cubic growth function

$$W(t) = C * (t - K)^3 + W_{max}$$

Di cui:

- C è un parametro che determina l'aggressività dell'algoritmo, sia in fase di incremento verso il punto di saturazione, sia in fase di probing.
- t rappresenta il tempo trascorso dall'ultima entrata in *Fast Recovery*, in cui viene fissato il valore per W_{max} .
- K rappresenta il tempo necessario per portare la finestra di congestione da W a W_{max} .

Sono piccole le differenze rispetto la Bic function, in quanto è sicuramente importante mantenere il buon comportamento di Bic sulle reti caratterizzate da un elevato *BDP*.

Vi è tuttavia un'altra differenza, attraverso la quale il Cubic riesce a comportarsi meglio, e ad essere più fair verso i flussi Reno, in condizioni di basso RTT, ossia la *Tcp Mode*. Infatti nel corso della fase di *Congestion Avoidance*, Cubic incrementa ad ogni RTT la *cwnd* sulla base della cubic function, solo se la *cwnd* si trova al di sopra del valore che si sarebbe ottenuto con l'approccio standard (s'intende Reno, new Reno).

In caso contrario si è evidentemente in una condizione in cui l'incremento lineare per ogni RTT, funziona meglio. Per cui il Cubic entra in *Tcp Mode*, ossia opera come l'approccio standard.

Bibliografia

- [Abr15] Mikael Abrahamsson. *TCP ACK Suppression*. IETF AQM mailing list. Nov. 2015. URL: <https://www.ietf.org/mail-archive/web/aqm/current/msg01480.html>.
- [Car+17a] Neal Cardwell et al. *BBR Congestion Control: An Update*. Slides. Mar. 2017. URL: <https://www.ietf.org/proceedings/98/slides/slides-98-iccr-g-an-update-on-bbr-congestion-control-00.pdf>.
- [Car+17b] Neal Cardwell et al. «BBR: Congestion-based Congestion Control». In: *Commun. ACM* 60.2 (gen. 2017), pp. 58–66. ISSN: 0001-0782. DOI: 10.1145/3009824. URL: <http://doi.acm.org/10.1145/3009824>.
- [HRX08] Sangtae Ha, Injong Rhee e Lisong Xu. «CUBIC: A New TCP-friendly High-speed TCP Variant». In: *SIGOPS Oper. Syst. Rev.* 42.5 (lug. 2008), pp. 64–74. ISSN: 0163-5980. DOI: 10.1145/1400097.1400105. URL: <http://doi.acm.org/10.1145/1400097.1400105>.
- [RSc17] I.Rhee NCSU L.Xu UNL S.Ha Colorado A.Zimmermann L.Eggert NetApp R.Scheffenegger. *CUBIC for Fast Long-Distance Networks*. Internet-Draft draft-ietf-tcpm-cubic-05.txt. IETF Secretariat, lug. 2017.
- [YCh17a] Inc Y.Cheng N.Cardwell S.Hassas Yeganeh V.Jacobson Google. *BBR Congestion Control*. Internet-Draft draft-cardwell-iccr-g-bbr-congestion-control-00.txt. IETF Secretariat, lug. 2017.
- [YCh17b] Inc Y.Cheng N.Cardwell S.Hassas Yeganeh V.Jacobson Google. *Delivery Rate Estimation*. Internet-Draft draft-cheng-iccr-g-delivery-rate-estimation-00.txt. IETF Secretariat, lug. 2017.