WebRTC基于GCC的拥塞控制(上) - 算法分析

jianshu.com/p/0f7ee0e0b3be

实时流媒体应用的最大特点是实时性,而延迟是实时性的最大敌人。从媒体收发端来讲,媒体数据的处理速度是造成延迟的重要原因;而从传输角度来讲,网络拥塞则是造成延迟的最主要原因。网络拥塞可能造成数据包丢失,也可能造成数据传输时间变长,延迟增大。

拥塞控制是实时流媒体应用质量保证(QoS)的重要手段之一,它在缓解网络拥堵、减小网络延迟、平滑数据传输等质量保证方面发挥重要作用。WebRTC通控制发送端数据发送码率来达到控制网络拥塞的目的,其采用谷歌提出的拥塞控制算法(Google Congestion Control,简称GCC[1])来控制发送端码率。

本文是关于WebRTC拥塞控制算法GCC的上半部分,主要集中于对算法的理论分析,力图对WebRTC的QoS有一个全面直观的认识。在下半部分,将深入WebRTC源代码内部,仔细分析GCC的实现细节。

1 GCC算法综述

Google关于GCC的RFC文档在文献[1],该RFC目前处于草案状态,还没有成为IETF的正式RFC。此外,Google 陆续发布了一系列论文[2][3][4]来论述该算法的实现细节,以及其在Google Hangouts、WebRTC等产品中的应用。本文主要根据这些文档资料,从理论上学习GCC算法。

GCC算法分两部分:发送端基于丢包率的码率控制和接收端基于延迟的码率控制。如图1所示。

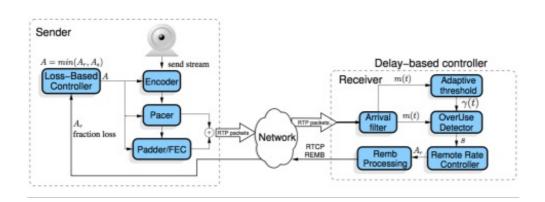


图1 GCC算法整体结构

基于丢包率的码率控制运行在发送端,依靠RTCP RR报文进行工作。WebRTC在发送端收到来自接收端的RTCP RR报文,根据其Report Block中携带的丢包率信息,动态调整发送端码率As。基于延迟的码率控制运行在接收端,WebRTC根据数据包到达的时间延迟,通过到达时间滤波器,估算出网络延迟m(t),然后经过过载检测器判断当前网络的拥塞状况,最后在码率控制器根据规则计算出远端估计最大码率Ar。得到Ar之后,通过RTCP REMB报文返回发送端。发送端综合As、Ar和预配置的上下限,计算出最终的目标码率A,该码率会作用到 Encoder、RTP和PacedSender等模块,控制发送端的码率。

2 发送端基于丢包率的码率控制

GCC算法在发送端基于丢包率控制发送码率,其基本思想是:丢包率反映网络拥塞状况。如果丢包率很小或者为 0,说明网络状况良好,在不超过预设最大码率的情况下,可以增大发送端码率;反之如果丢包率变大,说明网络 状况变差,此时应减少发送端码率。在其它情况下,发送端码率保持不变。

GCC使用的丢包率根据接收端RTP接收统计信息计算得到,通过RTCP RR报文中返回给发送端。RTCP RR报文统计接收端RTP接收信息,如Packet Loss,Jitter,DLSR等等,如图2所示:

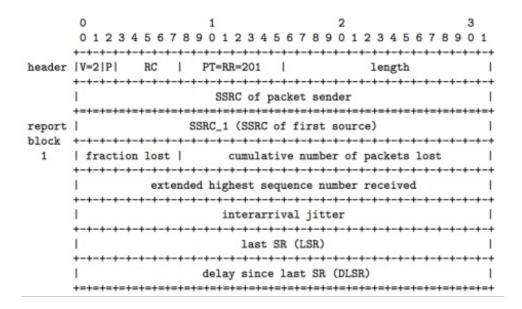


图2 RTCP RR报文结构[5]

发送端收到RTCP RR报文并解析得到丢包率后,根据图3公式计算发送端码率:当丢包率大于0.1时,说明网络发生拥塞,此时降低发送端码率;当丢包率小于0.02时,说明网络状况良好,此时增大发送端码率;其他情况下,发送端码率保持不变。

$$A_s(t_k) = \begin{cases} A_s(t_{k-1})(1 - 0.5f_l(t_k)) & f_l(t_k) > 0.1 \\ 1.05(A_s(t_{k-1})) & f_l(t_k) < 0.02 \\ A_s(t_{k-1}) & \text{otherwise} \end{cases}$$

图3 GCC基于丢包率的码率计算公式[4]

最终码率会作用于Encoder、RTP和PacedSender模块,用以在编码器内部调整码率和平滑发送端发送速率。

3 接收端基于延迟的码率控制

GCC算法在接收端基于数据包到达延迟估计发送码率Ar,然后通过RTCP REMB报文反馈到发送端,发送端把Ar作为最终目标码率的上限值。其基本思想是: RTP数据包的到达时间延迟m(i)反映网络拥塞状况。当延迟很小时,说明网络拥塞不严重,可以适当增大目标码率;当延迟变大时,说明网络拥塞变严重,需要减小目标码率;当延迟维持在一个低水平时,目标码率维持不变。

基于延时的拥塞控制由三个主要模块组成:到达时间滤波器,过载检查器和速率控制器;除此之外还有过载阈值 自适应模块和REMB报文生成模块,如图1所示。下面分别论述其工作过程。

3.1 到达时间滤波器(Arrival-time Filter)

该模块用以计算相邻相邻两个数据包组的网络排队延迟m(i)。数据包组定义为一段时间内连续发送的数据包的集合。一系列数据包短时间里连续发送,这段时间称为突发时间,建议突发时间为5ms。不建议在突发时间内的包间隔时间做度量,而是把它们做为一组来测量。通过相邻两个数据包组的发送时间和到达时间,计算得到组间延迟d (i)。组间延迟示意图及计算公式如图4所示:

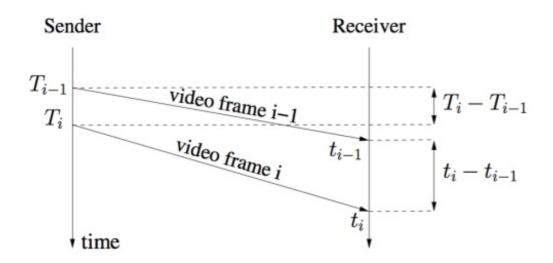


图4组间延迟示意图

T(i)是第i个数据包组中第一个数据包的发送时间,t(i)是第i个数据包组中最后一个数据包的到达时间。帧间延迟通过如下公式计算得到:

$$d(i) = t(i) - t(i-1) - (T(i) - T(i-1))$$
 (3.1.1)

公式1.3.1是d(i)的观测方程。另一方面,d(i)也可由如下状态方程得到:

$$d(i) = dL(i)/C(i) + w(i)$$

$$d(i) = dL(i)/C(i) + m(i) + v(i)$$
(3.1.2)
(3.1.3)

其中dL(i)表示相邻两帧的长度差,C(i)表示网络信道容量,m(i)表示网络排队延迟,v(i)表示零均值噪声。m(i)即是我们要求得的网络排队延迟。通过Kalman Filter可以求得该值。具体计算过程请参考文献[1][4][6]。

3.2 过载检测器(Over-use Detector)

该模块以到达时间滤波器计算得到的网络排队延迟m(i)为输入,结合当前阈值gamma_1,判断当前网络是否过载。判断算法如图5所示[2]。

```
1 if |m_i| > \gamma then
           if m_i > 0 then
 2
                 t_{OU} \leftarrow t_{OU} + \Delta T;
 3
                 if t_{OU} > \bar{t}_{OU} then
 4
                        if m_i \geq m_{i-1} then
 5
                              t_{OU} \leftarrow 0;
 s \leftarrow \text{Overuse};
 6
 7
           else
 8
                 t_{OU} \leftarrow 0;
 9
                  s \leftarrow \text{Underuse}:
10
11 else
           t_{OU} \leftarrow 0;
12
           s \leftarrow \text{Normal};
13
```

图5 过载检测器伪代码

算法基于当前网络排队延迟m(i)和当前阈值gamma_1判断当前网络拥塞状况[2]:当m(i) > gamma_1时,算法计算处于当前状态的持续时间t(ou) = t(ou) + delta(t),如果t(ou)大于设定阈值gamma_2(实际计算中设置为10ms),并且m(i) > m(i-1),则发出网络过载信号Overuse,同时重置t(ou)。如果m(i)小于m(i-1),即使高于阀值gamma_1也不需要发出过载信号。当m(i) < -gamma_1时,算法认为当前网络处于空闲状态,发出网络低载信号Underuse。当 – gamma_1 <= m(i) <= gamma_1是,算法认为当前网络使用率适中,发出保持信号Hold。算法随着时间轴的计算过程可从图6中看到。

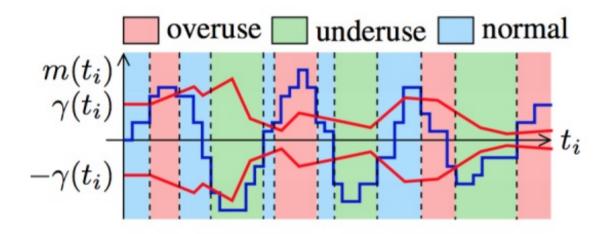


图6 时间轴上的过载检测过程

需要注意的是,阀值gamma_1对算法的影响很大,并且阈值gamma_1是自适应性的。如果其是静态值,会带来一系列问题,详见文献[4]。所以gamma_1需要动态调整来达到良好的表现。这就是图1中的Adaptive threshould 模块。阈值gamma_1动态更新的公式如下:

$$gamma_1(i) = gamma_1(i-1) + (t(i)-t(i-1)) * K(i) * (|m(i)|-gamma_1(i-1))$$
 (3.2.4)

当|m(i)|>gamma 1(i-1)时增加gamma 1(i),反之减小gamma 1(i),而当|m(i)|- gamma 1(i)>15,建议

gamma_1(i)不更新。K(i)为更新系数,当|m(i)|<gamma_1(i-1)时 $K(i) = K_d$,否则 $K(i) = K_u$ 。同时建议 gamma_1(i)控制在[6,600]区间。太小的值会导致探测器过于敏感。建议增加系数要大于减少系数 $K_u > K_d$ 。文献[1]给出的建议值如下:

gamma_1(0) = 12.5 ms gamma_2 = 10 ms K_u = 0.01 K_d = 0.00018

3.3 速率控制器(Remote Rate Controller)

该模块以过载检测器给出的当前网络状态s为输入,首先根据图7所示的有限状态机判断当前码率的变化趋势,然后根据图8所示的公式计算目标码率Ar。

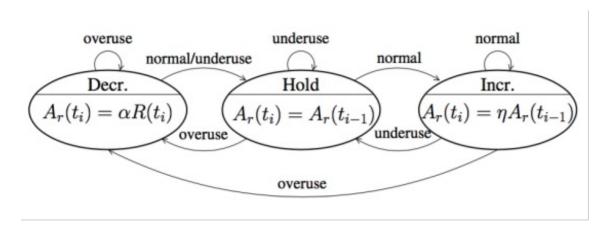


图7目标码率Ar变化趋势有限状态机

当前网络过载时,目标码率处于Decrease状态;当前网络低载时,目标码率处于Hold状态;当网络正常时,处于Decrease状态时迁移到Hold状态,处于Hold/Increase状态时都迁移到Increase状态。当判断出码率变化趋势后,根据图8所示公式进行计算目标码率。

$$A_r(t_i) = egin{cases} \eta A_r(t_{i-1}) & \sigma = ext{Increase} \\ \alpha R_r(t_i) & \sigma = ext{Decrease} \\ A_r(t_{i-1}) & \sigma = ext{Hold} \end{cases}$$

图8目标码率Ar计算公式

当码率变化趋势为Increase时,当前码率为上次码率乘上系数1.05;当码率变化趋势为Decrease,当前码率为过去500ms内的最大接收码率乘上系数0.85。当码率变化趋势为Hold时,当前码率保持不变。目标码率Ar计算得到之后,下一步把Ar封装到REMB报文中发送回发送端。REMB报文是Payload为206的RTCP报文[7],格式如图9所示。

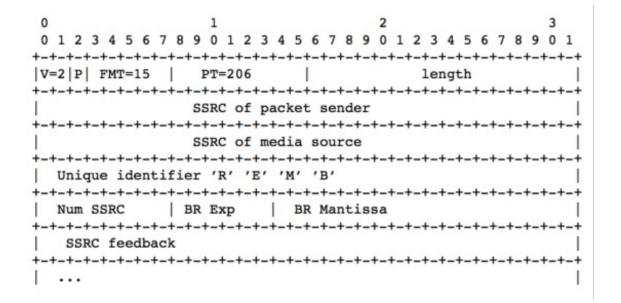


图9 REMB报文格式

REMB报文每秒发送一次, 当Ar(i) < 0.97 * Ar(i-1)时则立即发送。

3.4 发送端目标码率的确定

发送端最终目标码率的确定结合了基于丢包率计算得到的码率As和基于延迟计算得到的码率Ar。此外,在实际实现中还会配置目标码率的上限值和下限值。综合以上因素,最终目标码率确定如下:

```
target bitrate = max( min( min(As, Ar), Amax), Amin) (3.4.1)
```

目标码率确定之后,分别设置到Encoder模块和PacedSender模块。

4.总结

本文在广泛调研WebRTC GCC算法的相关RFC和论文的基础上,全面深入学习GCC算法的理论分析,以此为契机力图对WebRTC的QoS有一个全面直观的认识。为将来深入WebRTC源代码内部分析GCC的实现细节奠定基础。

参考文献

- [1] A Google Congestion Control Algorithm for Real-Time Communication.
- draft-alvestrand-rmcat-congestion-03
- [2] Understanding the Dynamic Behaviour of the Google Congestion Control for RTCWeb.
- [3] Experimental Investigation of the Google Congestion Control for Real-Time Flows.
- [4] Analysis and Design of the Google Congestion Control for Web Real-time Communication (WebRTC).

MMSys'16, May 10-13, 2016, Klagenfurt, Austria

- [5] RFC3550: RTP A Transport Protocol for Real-Time Applications
- [6] WebRTC视频接收缓冲区基于KalmanFilter的延迟模型.

http://www.jianshu.com/p/bb34995c549a

[7] RTCP message for Receiver Estimated Maximum Bitrate. draft-alvestrand-rmcat-remb-03