

# BALANCING TCP BUFFER SIZE VS PARALLEL STREAMS IN APPLICATION-LEVEL THROUGHPUT OPTIMIZATION

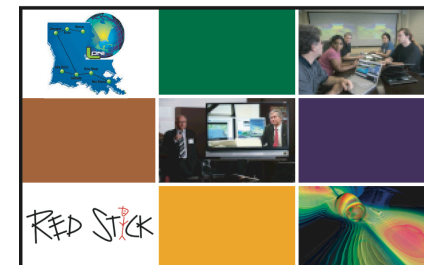
Esma Yildirim, Dengpan Yin, **Tevfik Kosar\***

Center for Computation & Technology  
Louisiana State University



CENTER FOR COMPUTATION  
& TECHNOLOGY

June 9, 2009  
DADC'09



# Motivation

- ▶ End-to-end data transfer performance is a major bottleneck for large-scale distributed applications
- ▶ TCP based solutions
  - Fast TCP, Scalable TCP etc
- ▶ UDP based solutions
  - RBUDP, UDT etc
- ▶ Most of these solutions require kernel level changes
- ▶ Not preferred by most domain scientists

# Application-Level Solution

- ▶ Take an application-level transfer protocol (i.e. GridFTP) and tune-up for optimal performance:
  - Using Multiple (Parallel) streams
  - Tuning Buffer size

# Roadmap

- ▶ Introduction
- ▶ Parallel Stream Optimization
- ▶ Buffer Size Optimization
- ▶ Combined Optimization of Buffer Size and Parallel Stream Number
- ▶ Conclusions

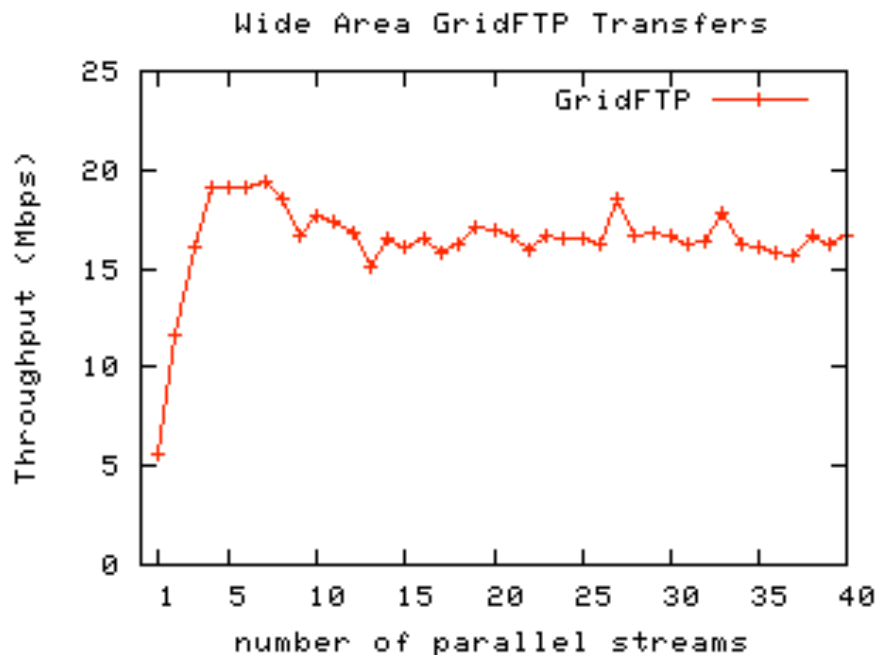


# Parallel Stream Optimization

**For a single stream**, theoretical calculation of throughput based on MSS, RTT and packet loss rate:

$$Th \leq \frac{MSS}{RTT} \frac{c}{\sqrt{p}}$$

**For n streams?**

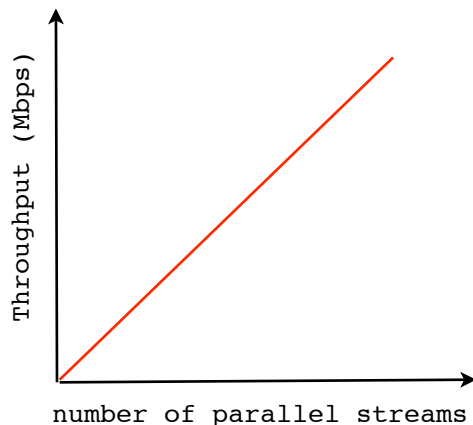


# Previous Models

## Hacker et al (2002)

An application opening  $n$  streams gains as much throughput as the total of  $n$  individual streams can get:

$$Th_n \leq \frac{MSS \times c}{RTT} \left( \frac{n}{\sqrt{p}} \right)$$

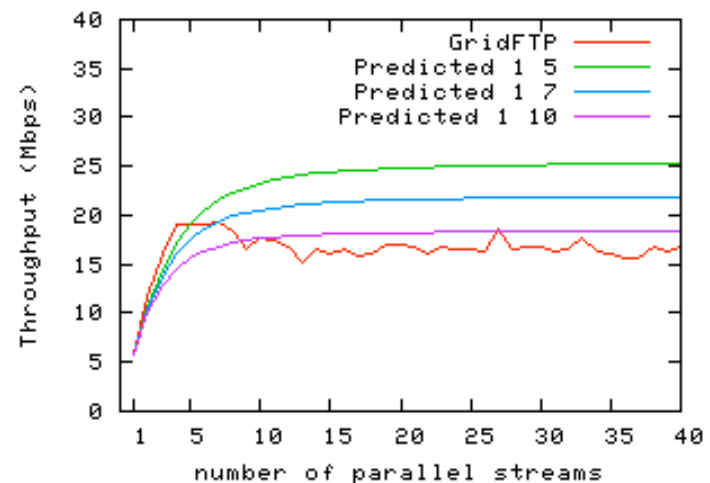


## Dinda et al (2005)

A relation is established between  $RTT$ ,  $p$  and the number of streams  $n$ :

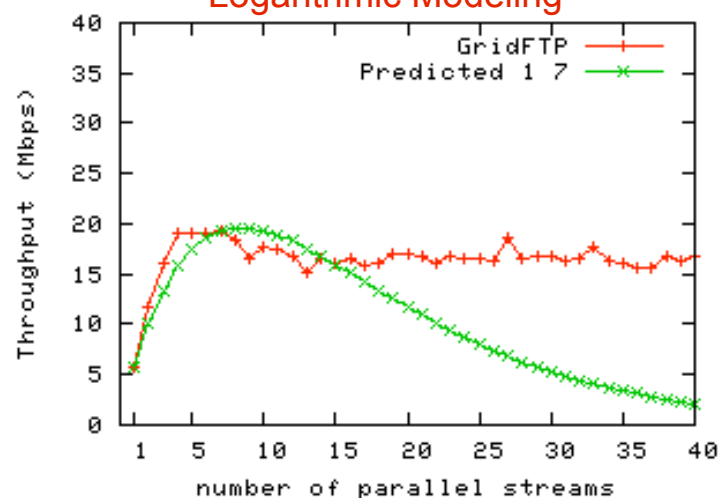
$$p'_n = p_n \frac{RTT_n^2}{c^2 MSS^2} = a'n^2 + b'$$

$$Th_n = \frac{n}{\sqrt{p'_n}} = \frac{n}{\sqrt{a'n^2 + b'}}$$

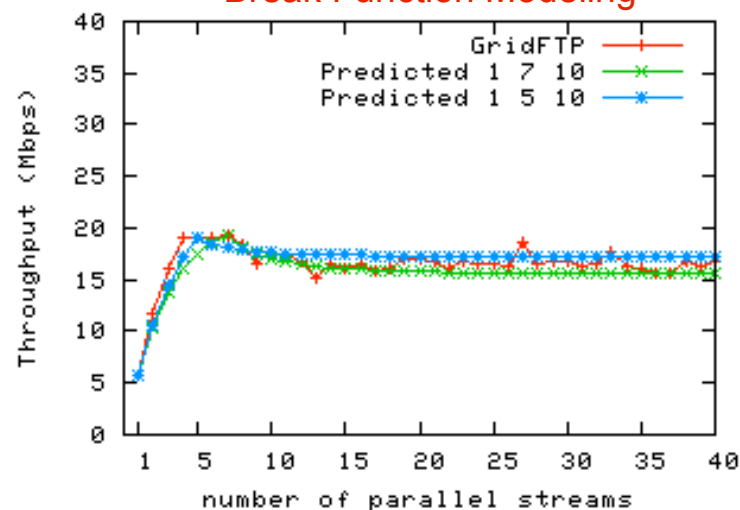


# Kosar et al Models

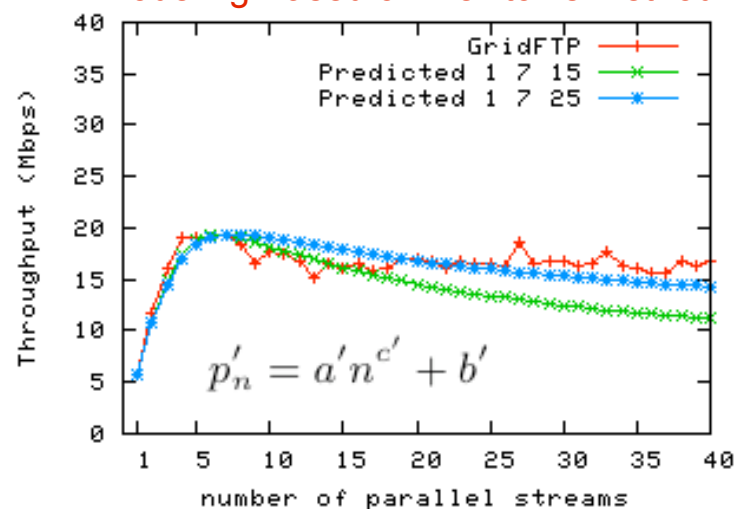
Logarithmic Modeling



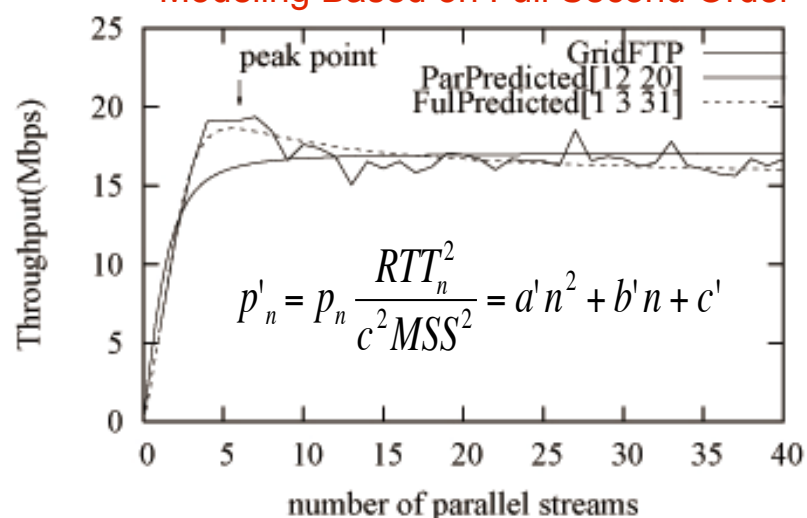
Break Function Modeling



Modeling Based on Newton's Method

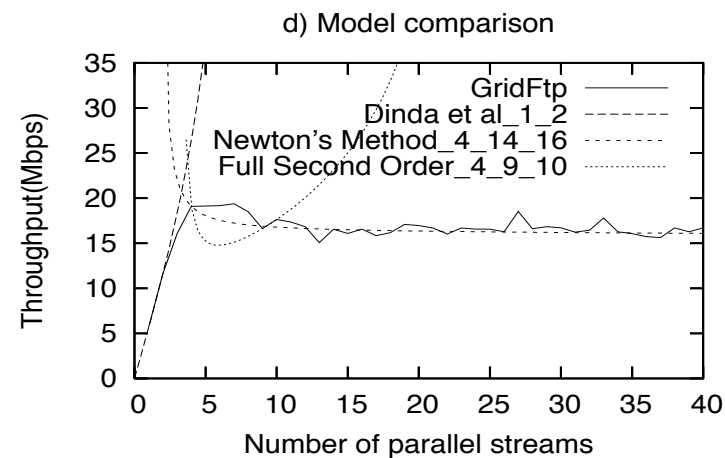
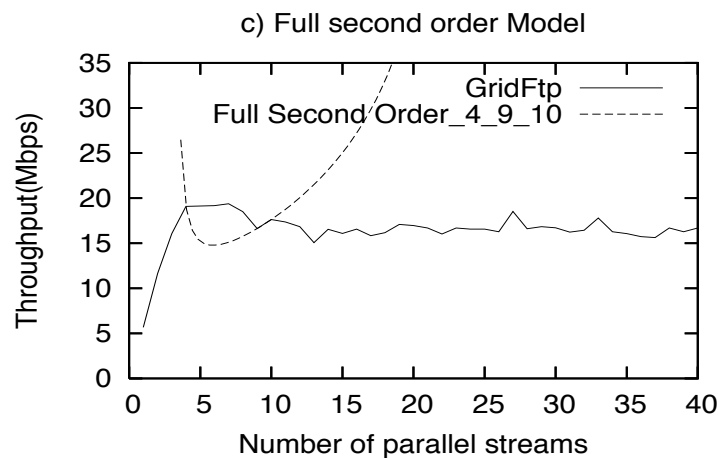
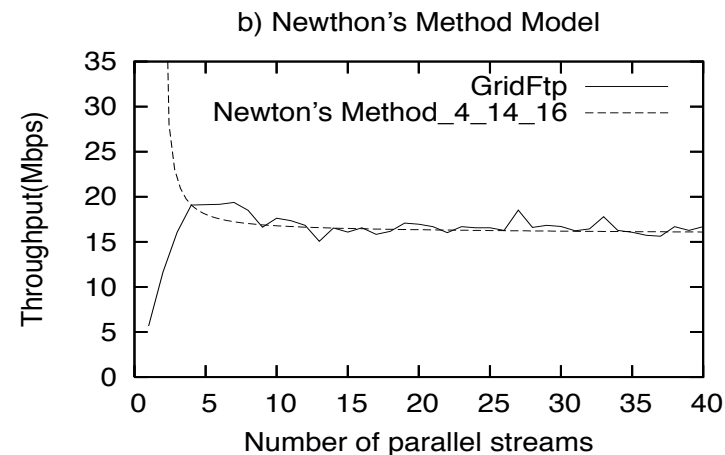
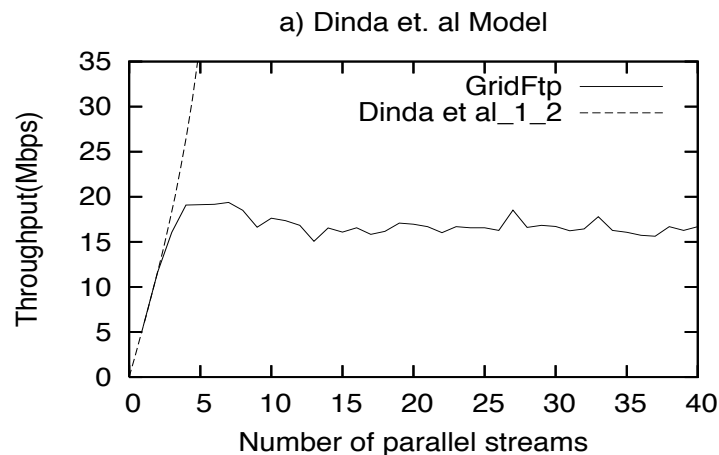


Modeling Based on Full Second Order



# It is not a perfect World!

- ▶ The selection of point should be made intelligently otherwise it could result in mispredictions



# Delimitation of Coefficients

- ▶ Pre-calculations of the coefficients of  $a'$ ,  $b'$  and  $c'$  and checking their ranges could save us for elimination of error rate

- ▶ Ex: Full second order

- $a' > 0$
- $b' < 0$
- $c' > 0$
- $2c' + b' > 1$

$$p'_n = p_n \frac{RTT_n^2}{c^2 MSS^2} = a' n^2 + b' n + c'$$

# Selection Algorithm

EXPSELECTION( $T$ )

```

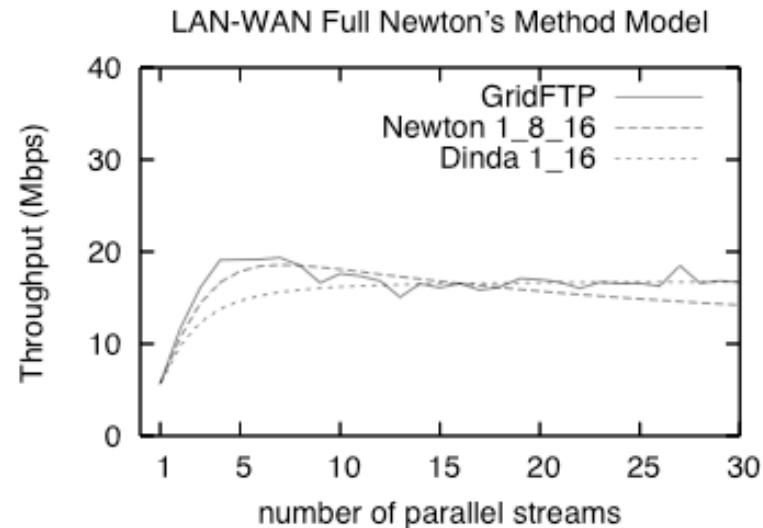
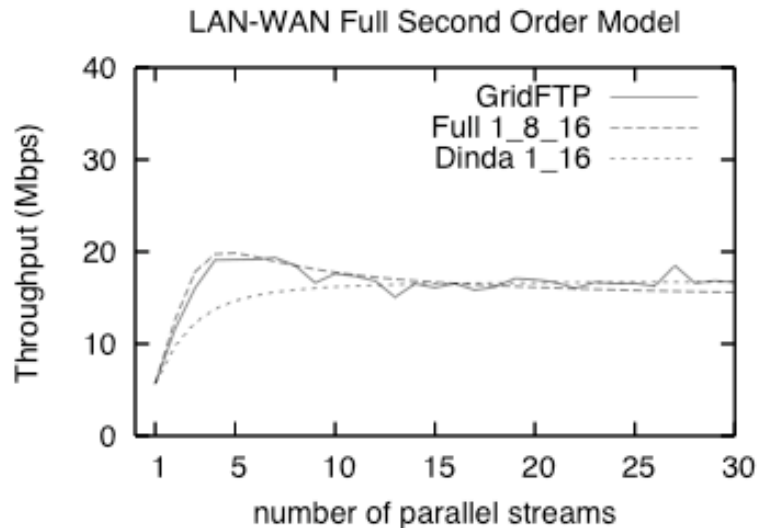
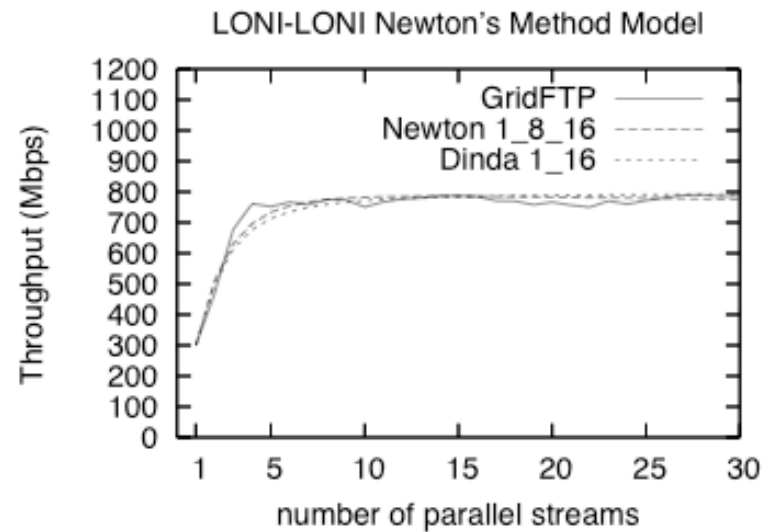
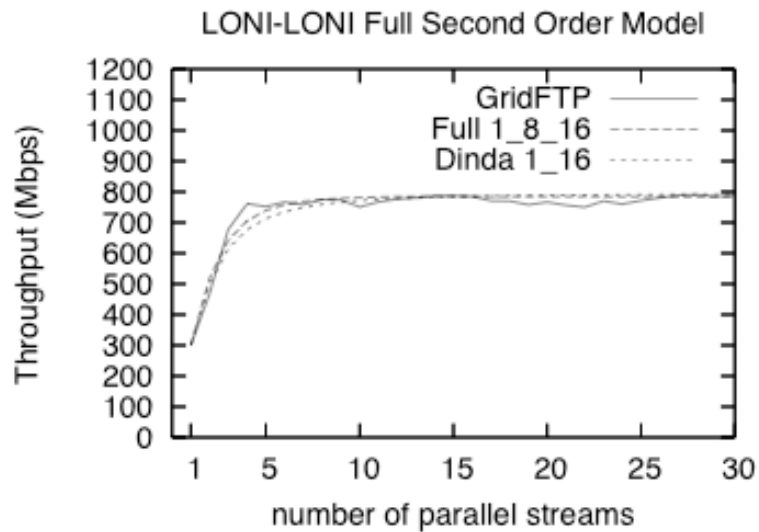
1  ▷ Input:  $T$ 
2  ▷ Output:  $O[i][j]$ 
3  1 Begin
4     $accuracy \leftarrow \alpha$ 
5     $i \leftarrow 1$ 
6     $streamno1 \leftarrow 1$ 
7     $throughput1 \leftarrow T_{streamno1}$ 
8     $O[i][1] \leftarrow streamno1$ 
9     $O[i][2] \leftarrow throughput1$ 
10   do
11      $streamno2 \leftarrow 2 * streamno1$ 
12      $throughput2 \leftarrow T_{streamno2}$ 
13      $slop \leftarrow \frac{throughput2 - throughput1}{streamno2 - streamno1}$ 
14      $i \leftarrow i + 1$ 
15      $O[i][1] \leftarrow streamno2$ 
16      $O[i][2] \leftarrow throughput2$ 
17      $streamno1 \leftarrow streamno2$ 
18      $throughput1 \leftarrow throughput2$ 
19   while  $slop > accuracy$ 
20 End
```

BESTCMB( $O, n, model$ )

```

1  ▷ Input:  $O, n$ 
2  ▷ Output:  $a, b, c, optnum$ 
3  1 Begin
4     $\overline{err}_m \leftarrow init$ 
5    for  $i \leftarrow 1$  to  $(n - 2)$  do
6      for  $j \leftarrow (i + 1)$  to  $(n - 1)$  do
7        for  $k \leftarrow (j + 1)$  to  $n$  do
8           $a', b', c' \leftarrow \text{CALCOE}(O, i, j, k, model)$ 
9          if  $a', b', c'$  are effective then
10              $\overline{err} \leftarrow \frac{1}{n} \sum_{t=1}^n |O[t][2] - Th_{pre}(O[t][1])|$ 
11             if  $\overline{err}_m = init \parallel err < err_m$  then
12                  $err_m \leftarrow err$ 
13                  $a \leftarrow a'$ 
14                  $b \leftarrow b'$ 
15                  $c \leftarrow c'$ 
16             end if
17         end if
18     end for
19 end for
20  $optnum \leftarrow \text{CALOPTSTREAMNO}(a, b, c, model)$ 
21 return  $optnum$ 
22 End
```

# Points Chosen by the Algorithm



# Buffer Size Optimization

- ▶ Buffer size affects the # of packets on the fly before an ack is received
- ▶ If undersized
  - The network can not be fully utilized
- ▶ If oversized
  - Throughput degradation due to packet losses which causes window reductions
- ▶ A common method is to set it to Bandwidth Delay Product = Bandwidth x RTT
- ▶ However there are differences in understanding the bandwidth and delay

# Bandwidth Delay Product

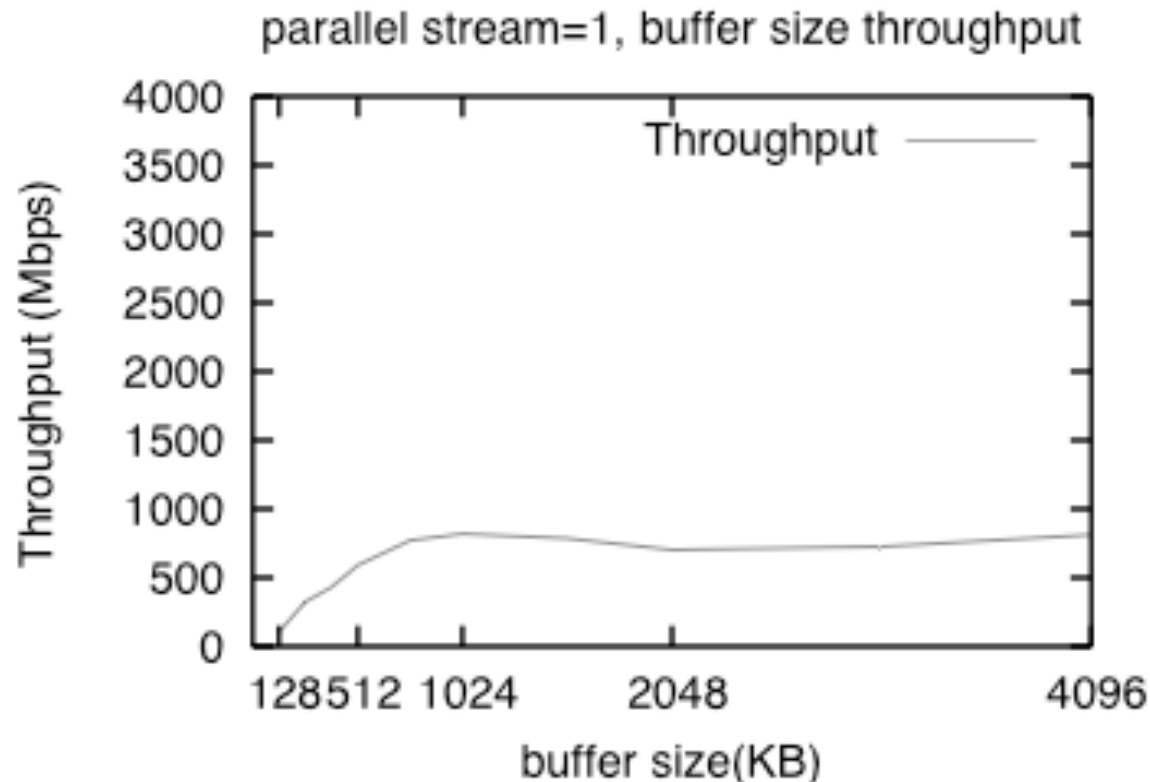
## ► BDP Types:

- $BDP1 = C \times RTT_{\max}$
- $BDP2 = C \times RTT_{\min}$ 
  - $C \rightarrow$  Capacity
- $BDP3 = A \times RTT_{\max}$
- $BDP4 = A \times RTT_{\min}$ 
  - $A \rightarrow$  Available bandwidth
- $BDP5 = BTC \times RTT_{\text{ave}}$ 
  - $BTC \rightarrow$  Average throughput of a congestion limited transfer
- $BDP6 = B_{\text{inf}}$ 
  - $B_{\text{inf}} \rightarrow$  a large value that is always greater than window size

# Existing Models

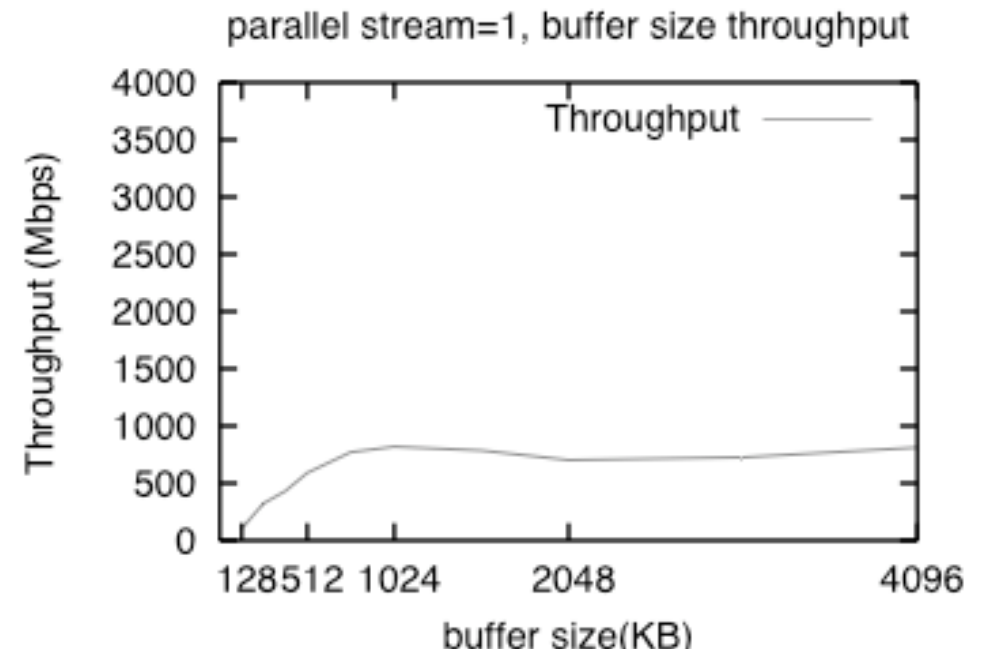
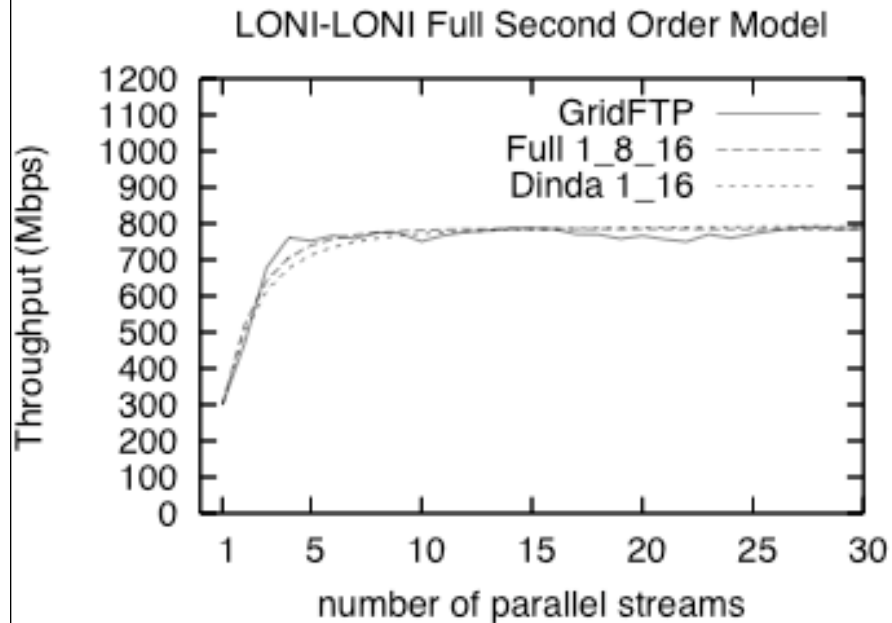
- ▶ Disadvantages of existing optimization techniques
  - Requires modification to the kernel
  - Rely on tools to take measurements of bandwidth and RTT
  - Do not consider the effect of cross traffic or congestion created by large buffer sizes
- ▶ Instead, can perform sampling and fit a curve to the buffer size graph

# Buffer Size Optimization



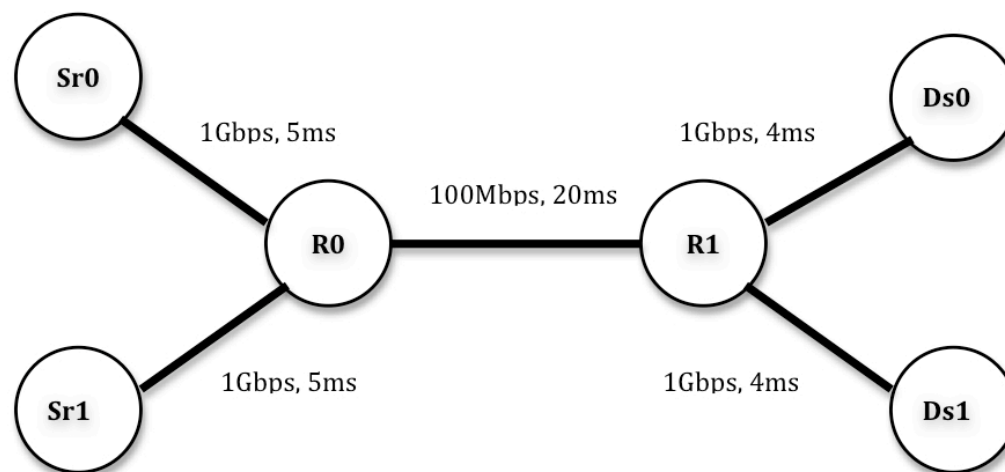
- ▶ Throughput becomes stable around 1M buffer size

# Combined Optimization



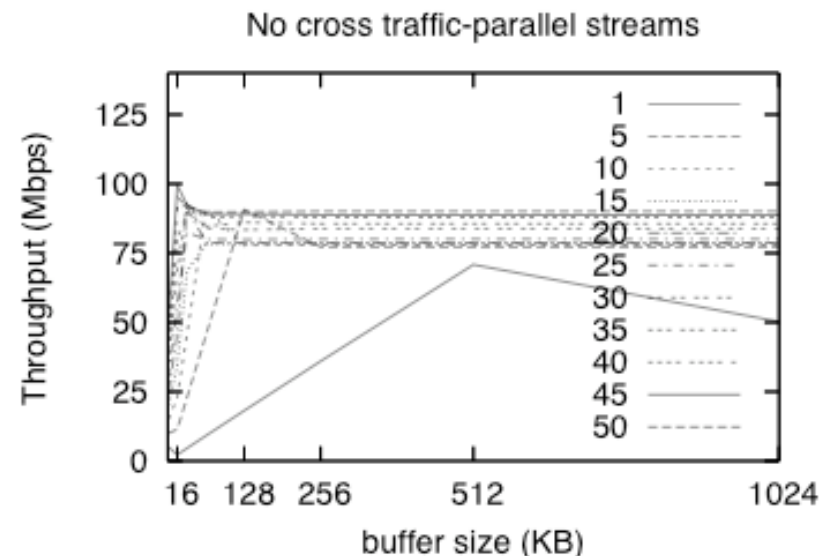
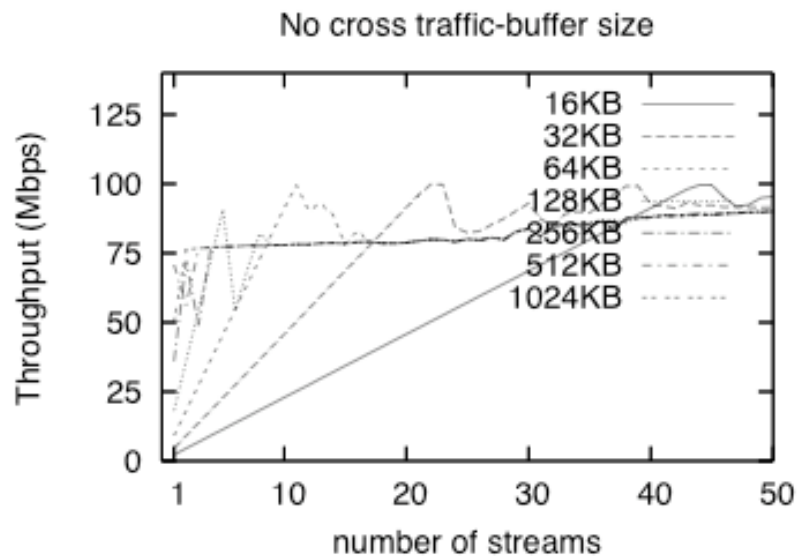
# Balancing: Simulations

- ▶ Simulator: NS-2
- ▶ Range of different buffer sizes and parallel streams used
- ▶ Test flows are from Sr1 to Ds1 where cross traffic is from Sr0 to Ds0



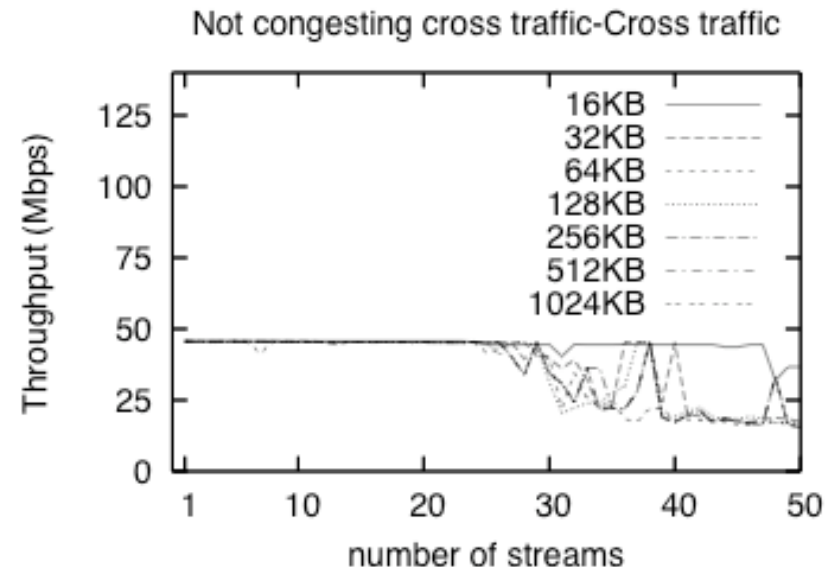
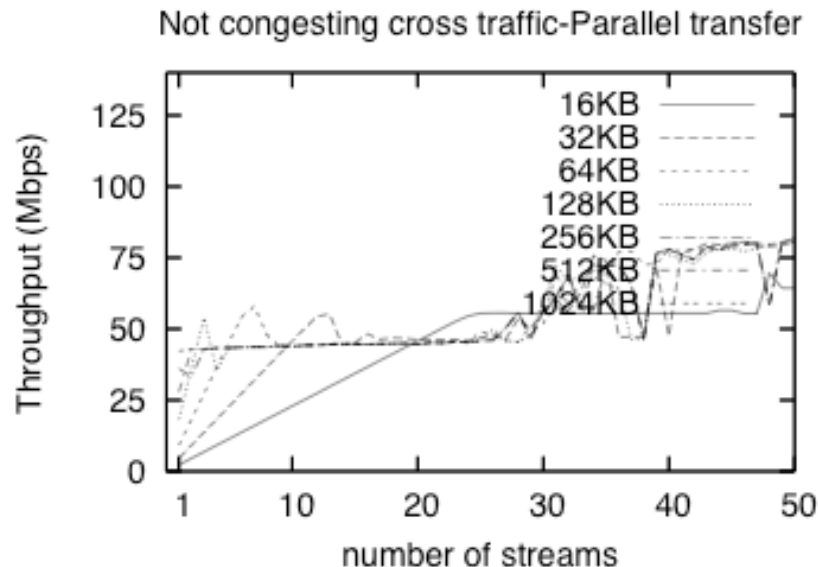
# 1 - No Cross Traffic

- ▶ Increasing the buffer size pulls back the parallel stream number to smaller values for peak throughput
- ▶ Further increasing the buffer size causes a drop in the peak throughput value



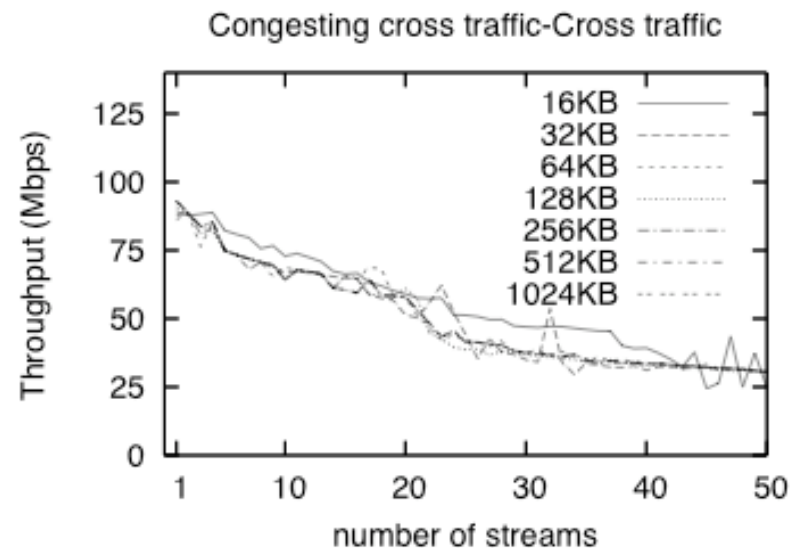
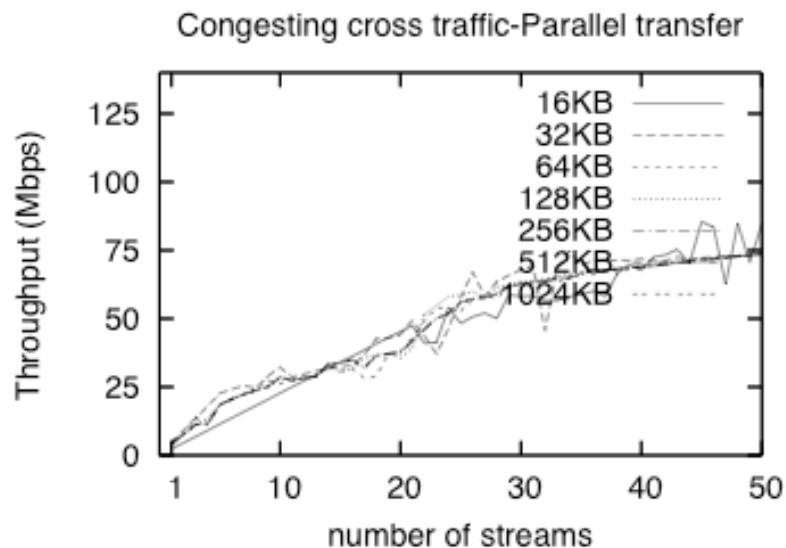
## 2 - Non-congesting Cross Traffic

- ▶ 5 streams of 64KB buffer size as traffic
- ▶ Similar behavior as no traffic case until the capacity is reached
- ▶ After the congestion starts the fight is won by the parallel flows of which stream number keeps increasing



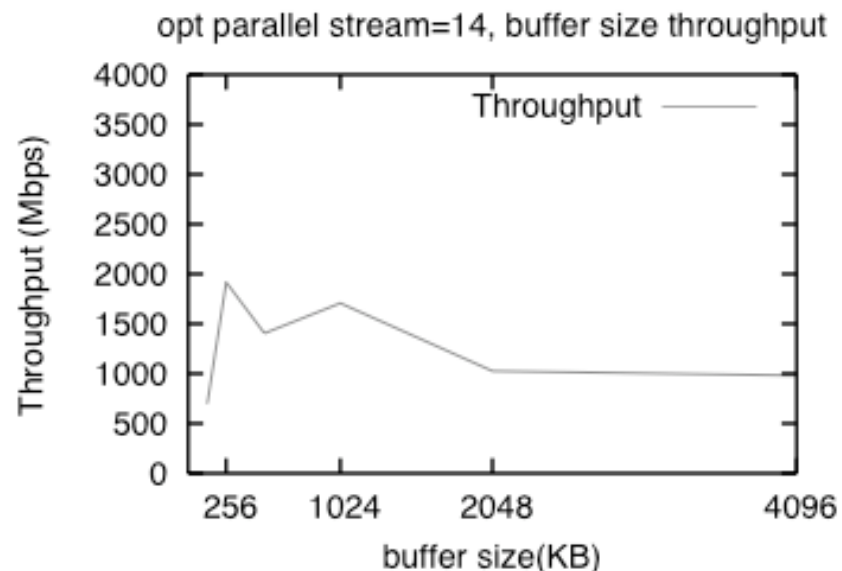
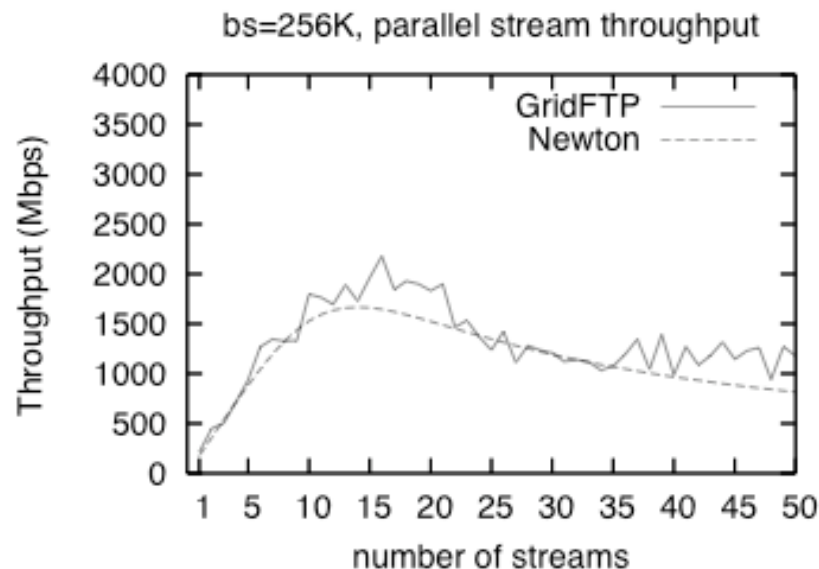
# 3 - Congesting Cross Traffic

- ▶ 12 streams of 64KB buffer size traffic
- ▶ No significant effect of buffer size
- ▶ As the number of parallel streams increases the throughput increases and cross traffic throughput decreases



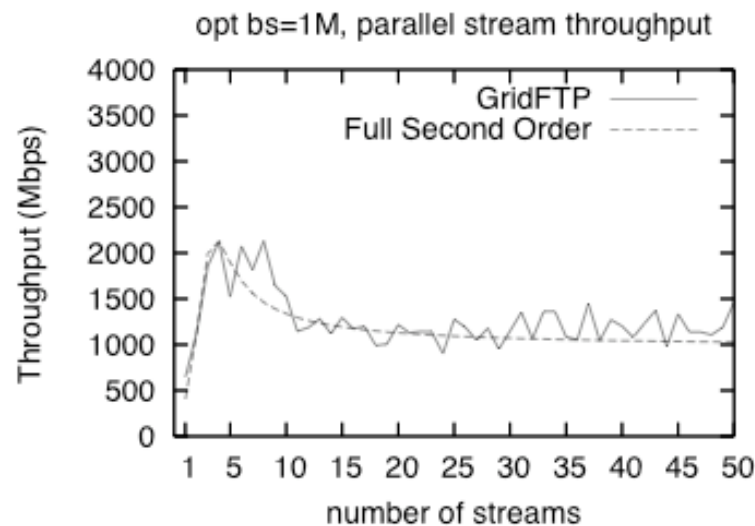
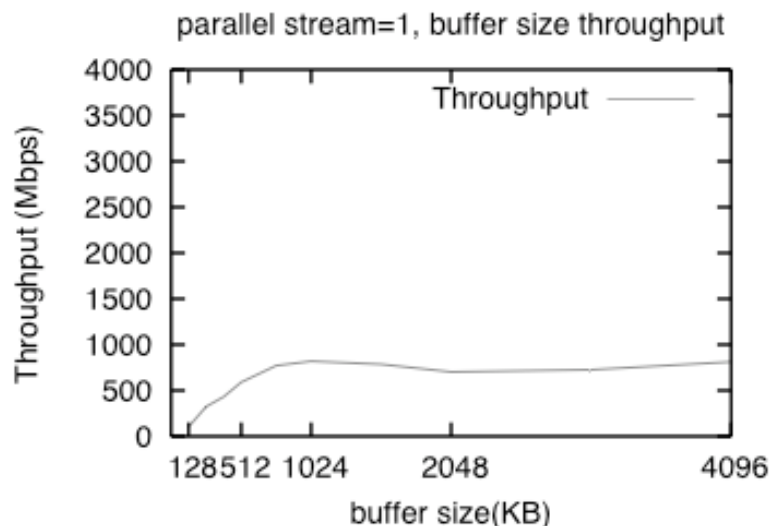
# Experiments on 10Gbps Network

- ▶ Approach 1: Tune # of streams first, then buffer size
  - Optimal stream number is 14 and an average peak of 1.7 Gbps is gained
  - Optimal buffer size = 256

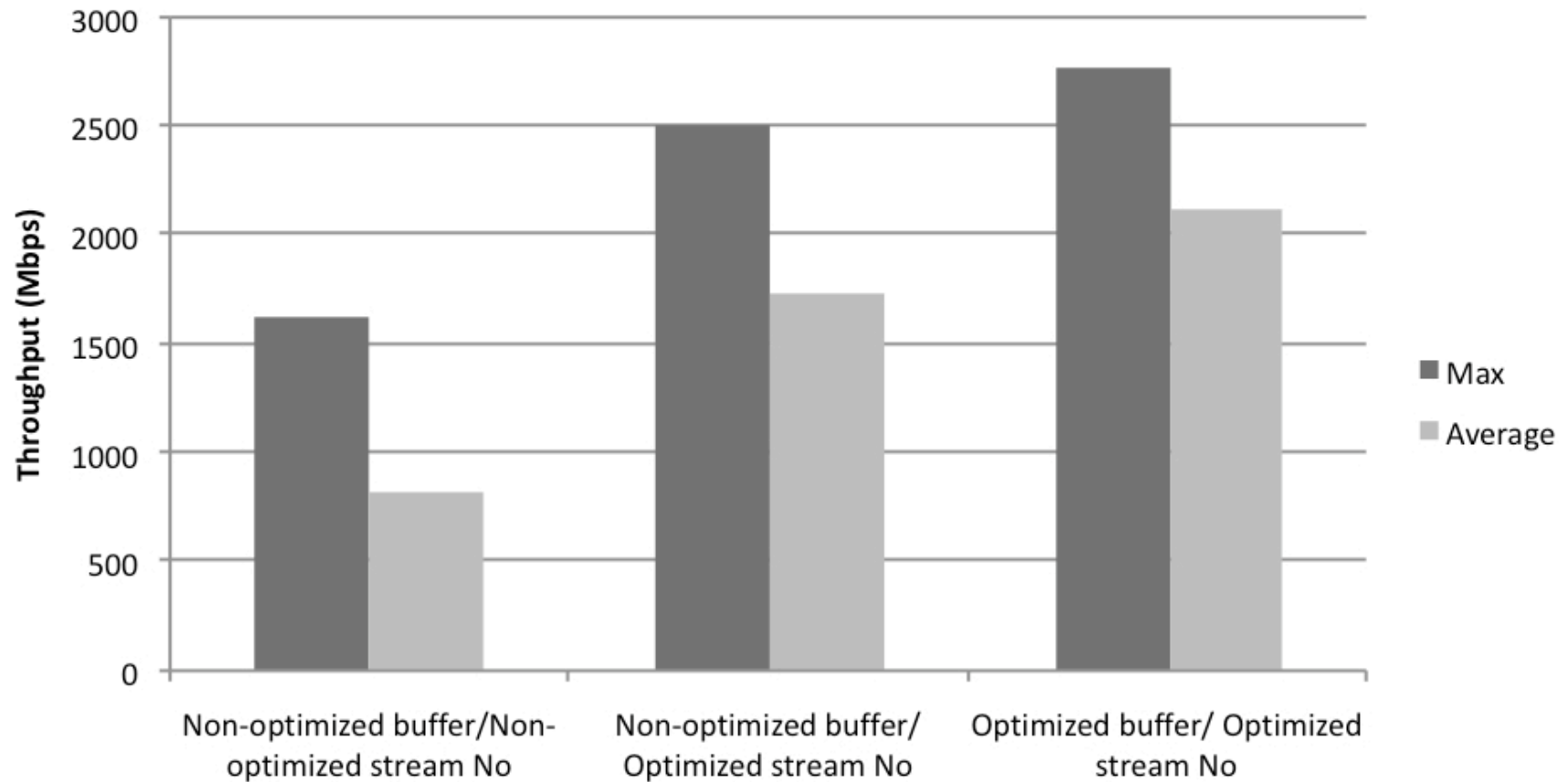


# Experiments on 10Gbps Network

- ▶ Approach 2: Tune buffer size first, then # of streams
  - Tuned buffer size for single stream is 1M and a throughput of around 900 Mbps is gained
  - Applying the parallel stream model, the optimal stream number is 4 and an average of around 2Gbps throughput is gained



## Comparison of Optimization Techniques



# Conclusions and Future Work

- ▶ Tuning buffer size and using parallel streams allow improvement of TCP throughput at the application level
- ▶ Two mathematical models (Newtons & Full Second Order) give promising results in predicting optimal number of parallel streams
- ▶ Early results in combined optimization show that using parallel streams on tuned buffers result in significant increase in throughput



This work has been sponsored by:  
**NSF and LA BoR**

For more information  
**Stork:** <http://www.storkproject.org>  
**PetaShare:** <http://www.petashare.org>