

# BwE: Flexible, Hierarchical Bandwidth Allocation for WAN Distributed Computing



ALOK KUMAR  
NIKHIL KASINADHUNI  
BJORN CARLIN

SUSHANT JAIN  
ENRIQUE ZERMENO  
MIHAI AMARANDEI-  
STAVILLA  
STEPHEN STUART

UDAY NAIK  
C. STEPHEN GUNN  
MATHIEU ROBIN  
AMIN VAHDAT

ANAND  
RAGHURAMAN  
JING AI  
ASPI SIGANPORIA

Google Inc.  
(Sigcomm'15)

25.05.2018

ETH Zurich

Presented by :  
Ritu Sriram

# Main idea



- **Service-level bandwidth allocation**
  - Prioritized bandwidth functions
  - Isolation between services
- **Suitable for distributed computing applications running across dedicated private WANs**
- **Mapping of bandwidth allocation policies onto packets at the host machines**

# Existing Approaches



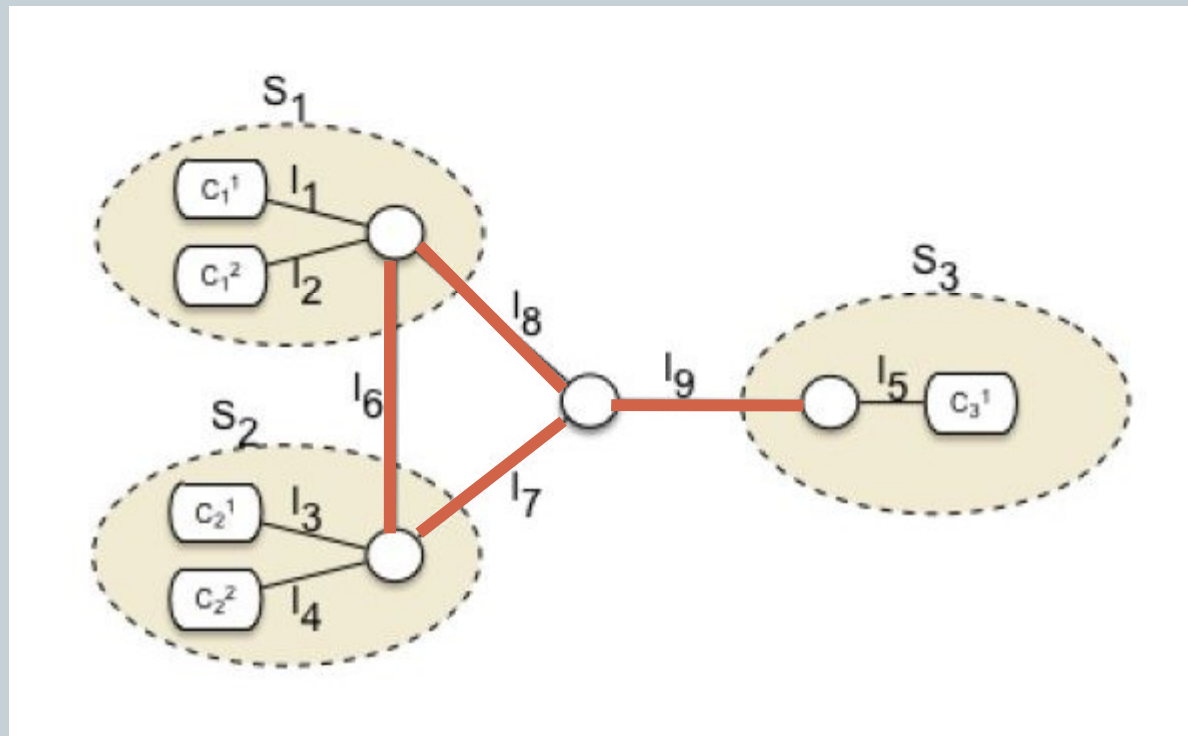
- **TCP-based bandwidth allocation**
  - All flows are of equal priority
  - Allocation is proportional to the # of active flows
- **QoS and MPLS tunnels**
  - Fixed the equal priority assumption
  - Scaling issues
  - Not enough flexibility in allocation policy

# Contributions



- Unified, hierarchical control plane for bandwidth management at the end-hosts
- Integrating BwE with existing TE solutions
- Hierarchical max-min fair bandwidth allocation algorithm

# Concepts : WAN Network Model



S – sites  
C – clusters  
l – links

# Concepts : FlowGroups or Traffic Aggregates



- Different granularities

Cluster fg

**src cluster | dst cluster | User\_aggregate**

(Input to BwE algorithm)

User fg

**src cluster | dst cluster | User\_name**



Job fg

**src cluster | dst cluster | User\_name | Job\_name**

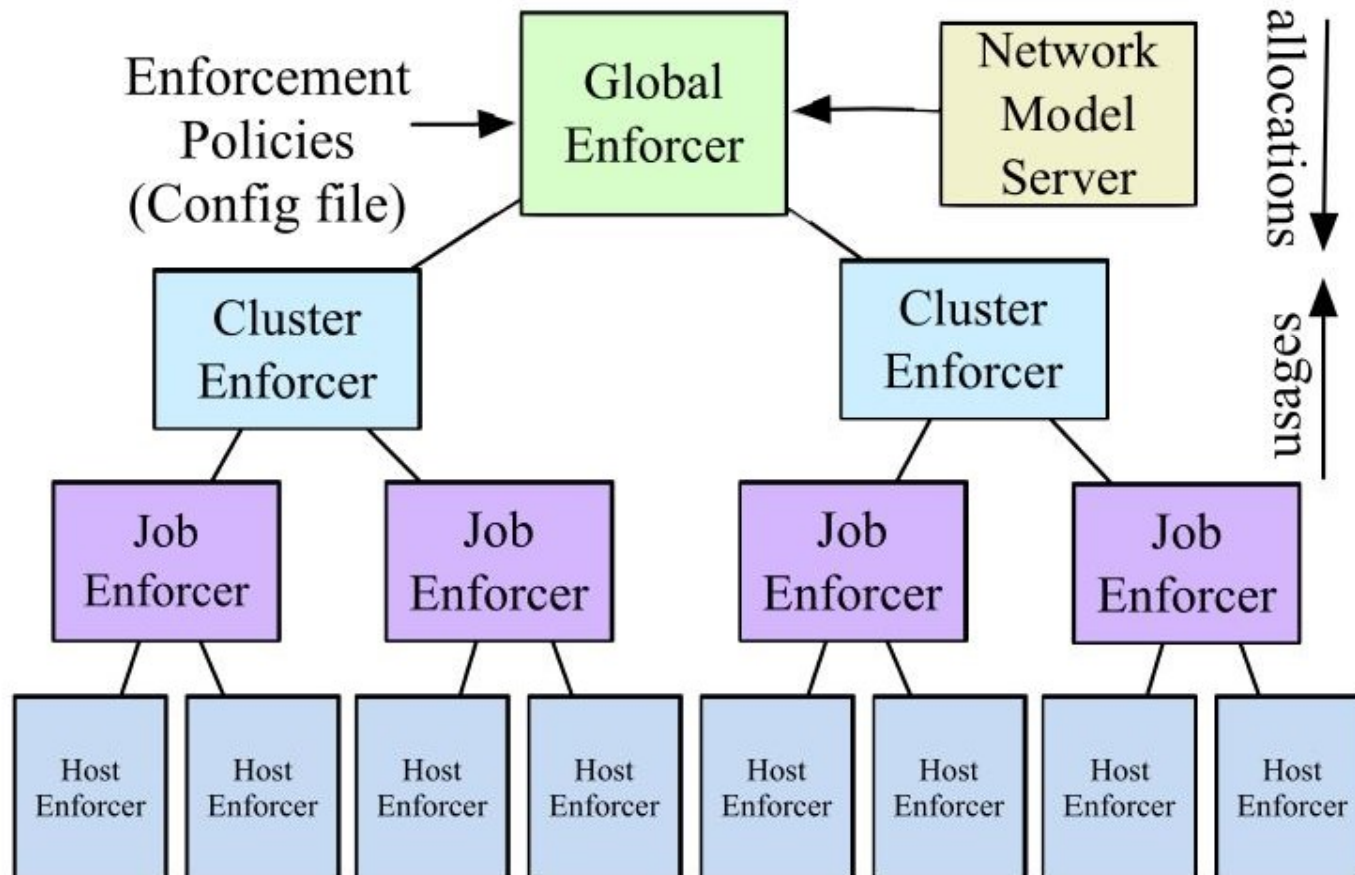


Task fg

**src cluster | dst cluster | User\_name | Job\_name | Task\_name**

(Measured and enforced at the hosts)

# BwE Architecture



# Concepts : BwE configuration



(a)  $f g_1$

Allocation Level	Weight	Bandwidth (Gbps)
Guaranteed	0	0
Best-Effort	20	10
	5	$\infty$

(b)  $f g_2$

Allocation Level	Weight	Bandwidth (Gbps)
Guaranteed	10	10
Best-Effort	10	$\infty$

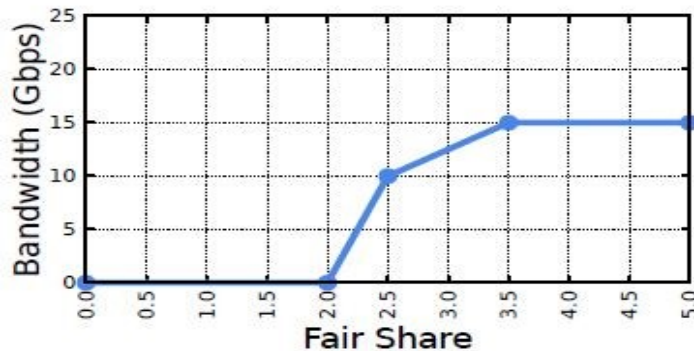


# Concepts : Bandwidth functions

- Specifies the bandwidth allocation to a FlowGroup as a function of its relative priority and 'fair-share'

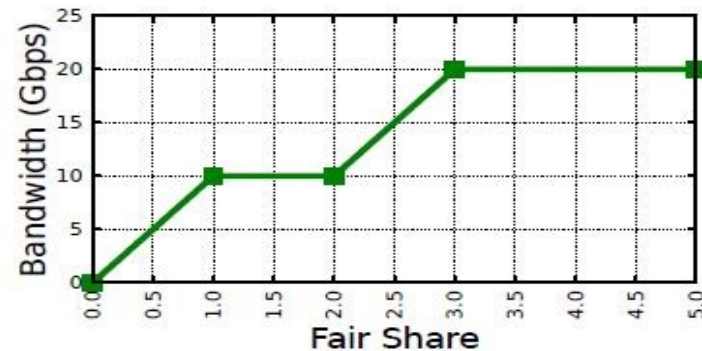
(a)  $f g_1$

Allocation Level	Weight	Bandwidth (Gbps)
Guaranteed	0	0
Best-Effort	20	10
	5	$\infty$



(b)  $f g_2$

Allocation Level	Weight	Bandwidth (Gbps)
Guaranteed	10	10
Best-Effort	10	$\infty$



Fair-share: measure of available capacity to be shared fairly

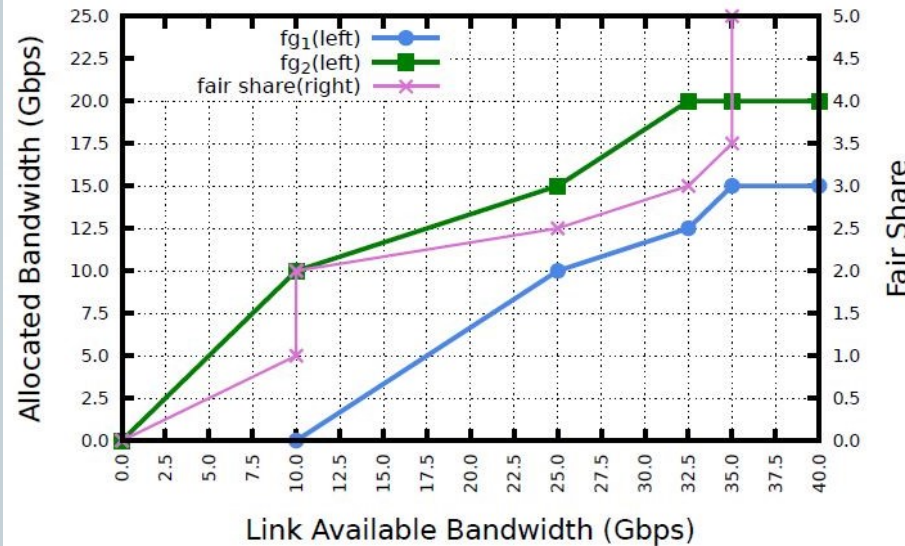
# BwE algorithm : Global MPFA



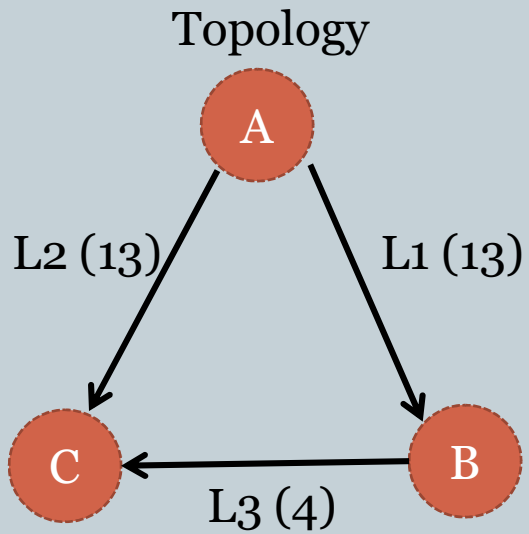
- **Inputs**

- Task fg's demands aggregated to the top level
- Network topology , optimal paths and weights from the underlying TE system

- **Output**



# MPFA Problem : Setup



## Inputs

Flow Group	Src	Dst	Paths
$f_1$	A	C	$l_2: 0.75$ $l_1 \rightarrow l_3: 0.25$
$f_2$	A	C	$l_2: 1.00$
$f_3$	A	B	$l_1: 1.00$

- All flows have demand – 18
- Weights of flows (priorities) :
  - $f_1$  – 1
  - $f_2$  – 2
  - $f_3$  – 3

\* Link capacities in Gbps

# MPFA : Bandwidth functions (FlowGroups)



- $B_{f_1}(s) = \min(18, s)$
- $B_{f_2}(s) = \min(18, 2s)$
- $B_{f_3}(s) = \min(18, 3s)$

Flow Group	Src	Dst	Paths
$f_1$	A	C	$l_2: 0.75$ $l_1 \rightarrow l_3: 0.25$
$f_2$	A	C	$l_2: 1.00$
$f_3$	A	B	$l_1: 1.00$

- All flows have demand – 18
- Weights of flows:
  - $f_1 - 1$
  - $f_2 - 2$
  - $f_3 - 3$

$s$  – fair share measure

# MPFA : Fraction of traverse



- Formula: Sum of weights for all paths of a fg containing the given link

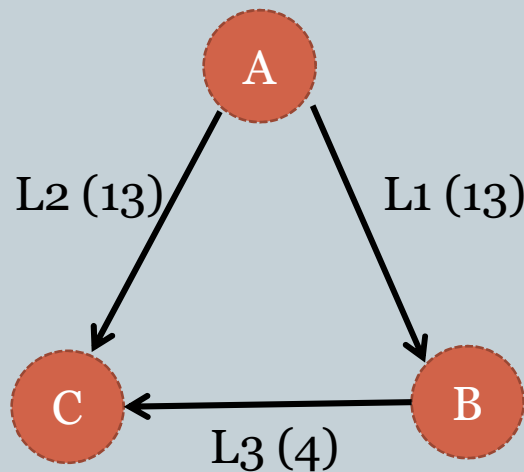
$$FR(f_i, l_k) = \sum_{1 \leq j \leq n_{f_i} | l_k \in p_j^{f_i}} w_j^{f_i}$$

- $FR(f_1, l_1) = 0.25$   
 $FR(f_1, l_2) = 0.75$   
 $FR(f_1, l_3) = 0.25$
- $FR(f_2, l_2) = 1$
- $FR(f_3, l_1) = 1$

Flow Group	Src	Dst	Paths
f <sub>1</sub>	A	C	l <sub>2</sub> : 0.75 l <sub>1</sub> →l <sub>3</sub> : 0.25
f <sub>2</sub>	A	C	l <sub>2</sub> : 1.00
f <sub>3</sub>	A	B	l <sub>1</sub> : 1.00

- All flows have demand – 18
- Weights of flows:
  - f<sub>1</sub> – 1
  - f<sub>2</sub> – 2
  - f<sub>3</sub> – 3

# MPFA : Bandwidth functions (Links)



Flow Group	Src	Dst	Paths
$f_1$	A	C	$l_2: 0.75$ $l_1 \rightarrow l_3: 0.25$
$f_2$	A	C	$l_2: 1.00$
$f_3$	A	B	$l_1: 1.00$

- Aim : Allocate bandwidth fairly without exceeding link capacity

$$B_{f_1}(s) = \min(18, s) ; B_{f_3}(s) = \min(18, 3s) ;$$

$$FR(f_1, l_1) = 0.25 ; FR(f_3, l_1) = 1 ;$$

$$B_{l_1}(s) = \left( \begin{array}{l} 0.25(\min(18, s)) \\ + \min(18, 3s) \end{array} \right) = \begin{cases} 3.25s & : 0 \leq s < 6 \\ 0.25s + 18 & : 6 \leq s < 18 \\ 22.5 & : s \geq 18 \end{cases}$$

# MPFA : Bandwidth functions (Links)



$$B_{l_1}(s) = \left( \begin{array}{l} 0.25(\min(18, s)) \\ + \min(18, 3s) \end{array} \right) = \begin{cases} 3.25s & : 0 \leq s < 6 \\ 0.25s + 18 & : 6 \leq s < 18 \\ 22.5 & : s \geq 18 \end{cases}$$

$$B_{l_2}(s) = \left( \begin{array}{l} 0.75(\min(18, s)) \\ + \min(18, 2s) \end{array} \right) = \begin{cases} 2.75s & : 0 \leq s < 9 \\ 0.75s + 18 & : 9 \leq s < 18 \\ 31.5 & : s \geq 18 \end{cases}$$

$$B_{l_3}(s) = 0.25(\min(18, s)) = \begin{cases} 0.25s & : 0 \leq s < 18 \\ 4.5 & : s \geq 18 \end{cases}$$

# MPFA : Finding bottleneck



- We find a bottleneck fair share for each link such that its bandwidth reaches link capacity.

$$B_{l_1}(s) = \left( \begin{array}{l} 0.25(\min(18, s)) \\ + \min(18, 3s) \end{array} \right) = \begin{cases} 3.25s & : 0 \leq s < 6 \\ 0.25s + 18 & : 6 \leq s < 18 \\ 22.5 & : s \geq 18 \end{cases}$$

Link		Bottleneck fair share
1	$(3.25 * s) = 13$	$s = 4$ ;
2	$(2.75 * s) = 13$	$s = 4.72$ ;
3	$(0.25 * s) = 4$	$s = 16$ ;

L1 has reached **full capacity first**, it no longer participates in MPFA !

L1 is the bottleneck link



# MPFA : Update functions

Flow Group	Src	Dst	Paths
$f_1$	A	C	$l_2: 0.75$ $l_1 \rightarrow l_3: 0.25$
$f_2$	A	C	$l_2: 1.00$
$f_3$	A	B	$l_1: 1.00$

Frozen FlowGroups

- Need to update Bandwidth functions for links 2 and 3 beyond fair share of 4

Accounted for  $f_1$

We found  $s = 4$  when it has max fair allocation

$$B_{l_2}(s) = \left( \begin{array}{c} 0.75(\min(18, s)) \\ + \min(18, 2s) \end{array} \right)$$

# MPFA : Update functions

$$B_{l_2}(s) = \begin{cases} 2.75s & : 0 \leq s < 4 \\ 2s + 3 & : 4 \leq s < 9 \\ 21 & : s \geq 9 \end{cases}$$

$$B_{l_3}(s) = \begin{cases} 0.25s & : 0 \leq s < 4 \\ 1 & : s \geq 4 \end{cases}$$

Link		Bottleneck fair share
1	-	s = 4;
2	2s + 3 = 13	s = 5 ;
3	Never hits full capacity	s = ∞

L2 has reached **full capacity** next, it no longer participates in MPFA !  
L2 is the bottleneck link

# MPFA : Final output



- Terminate when all FlowGroups are frozen
- OUTPUT : allocation in fairshare = (f1:4, f2:5, f3:4)

$$B_{f_1}(s) = \min(18, 4) = 4 \text{ Gbps}$$

$$B_{f_2}(s) = \min(18, 2^*5) = 10 \text{ Gbps}$$

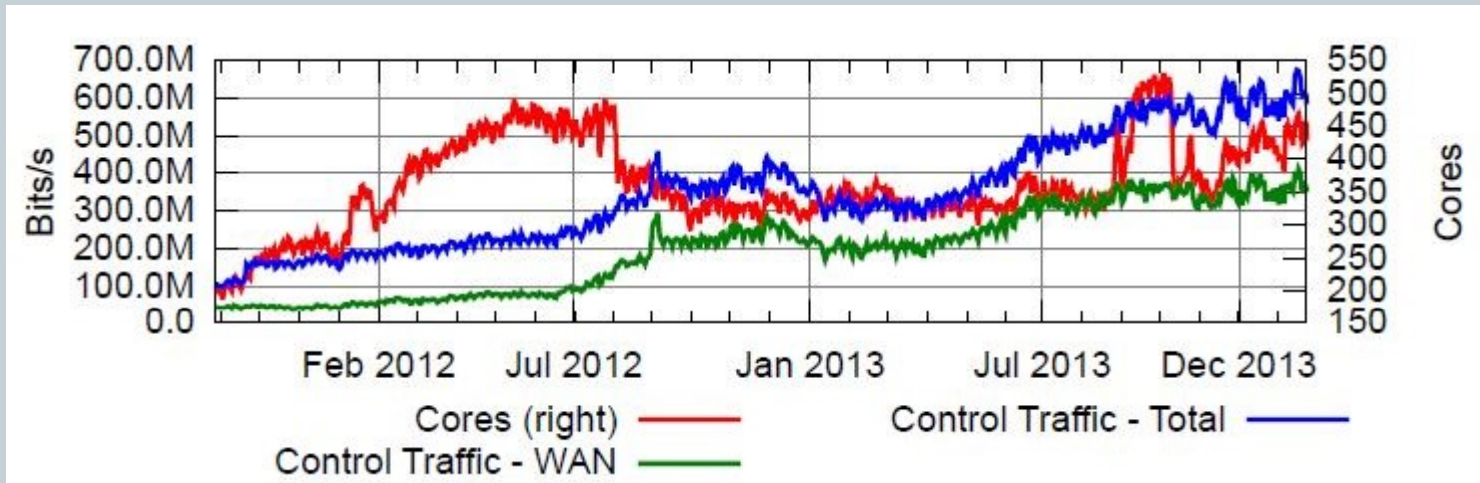
$$B_{f_3}(s) = \min(18, 3^*4) = 12 \text{ Gbps}$$

Best effort  
and  
As per priority

# Evaluation



- BwE has run in production for many years on Google WAN



- Algorithm run time

	Algo Run-time		Algo Interval(s)	Reporting Interval(s)
	Max(s)	Mean(s)		
Global Enforcer	3	-	10	10
Cluster Enforcer	.16	.15	4	10
Job Enforcer	<0.01	<0.01	4	5

# Conclusion



- Self-contained paper
- Well structured
- Mentions drawbacks and future optimization scope



Discussion ?