# Route Servers for !Dummies

## or: Scaling is Hard; Let's Go Shopping!

**inex**
*internet neutral exchange*

Nick Hilliard

Head of Operations

nick@inex.ie

# Some Blurb on INEX

- Currently only member-owner IXP in Ireland
- 59 members, 46 full members, 13 associate
- Estimate about 90% eyeballs in Ireland (South)
- Traffic levels: daytime peaks of 6G
- Provide usual services - 10M to 10G ethernet
- Two separate L2 infrastructures
- Three PoPs: Telecity Dublin, DEG, Interxion DUB1
- Mixture of Brocade (FES-X624, TI24X) and Cisco 6500
- Fibre lit with Transmode DWDM kit - N x 10G
- Highly active community interest

**i n e x**
internet  neutral  exchange

- Currently only member-owner IXP in Irela...
- 59 members, 46 full members, 1...
- Estimate about 90% ey...
- Traffic levels: d...
- Provid... ...G ethernet
- ...res
- ...ty Dublin, DEG, Interxion DUB1
- ...Brocade (FES-X624, TI24X) and Cisco 6500
- ...bre lit with Transmode DWDM kit - N x 10G
- Highly active community interest

*Free Beer at Meetings !!!11!!*

# Some Blurb on INEX

- Currently only member-owner IXP in Ireland
- 59 members, 46 full members, 13 associate
- Estimate about 90% eyeballs in Ireland (South)
- Traffic levels: daytime peaks of 6G
- Provide usual services - 10M to 10G ethernet
- Two separate L2 infrastructures
- Three PoPs: Telecity Dublin, DEG, Interxion DUB1
- Mixture of Brocade (FES-X624, TI24X) and Cisco 6500
- Fibre lit with Transmode DWDM kit - N x 10G
- Highly active community interest
- Oh yeah, we have some route-servers too

# Route Servers for Dummies

- Platform for multi lateral peering agreements (MPLA)
- Similar to a route reflector, except uses eBGP
- Very fashionable at IXPs right now
  - Reduce administrative load of peering
  - Simple interconnection to lots of other partners
  - Instant RoI (ISP management likes this)
  - Outsourcing RIB calculations to fast machines(!)
  - "Safe" if IXP has implemented prefix filtering
- Considered ghetto routing by larger providers
  - There are good reasons for this opinion
  - INEX recommends peering with route servers unless you know why you shouldn't
  - Because route servers are not for everyone
- Route prefix filtering considered indispensable by IXP participants
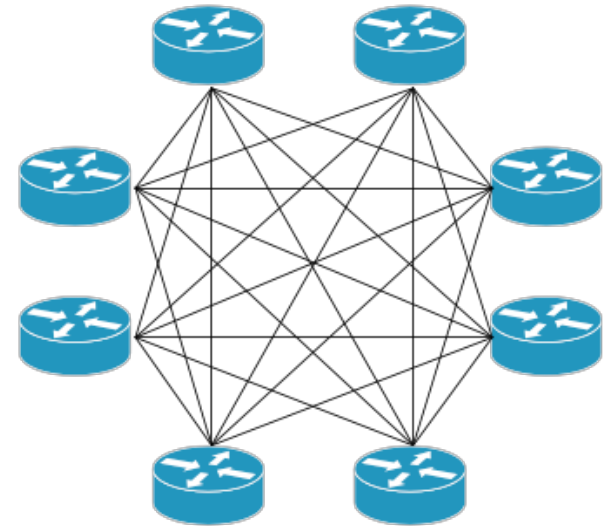
# Route Servers for Dummies

Peering on IXP without Route Servers

# Route Servers for Dummies

Peering on IXP without Route Servers

Peering on IXP with Route Servers

Route Server Cluster

# Single-RIB BGP policy problem

IXP Fabric

Y

A

X

Single RIB Route Server

C

B

D

☑ ABCD route server clients

# Single-RIB BGP policy problem

IXP Fabric

Y

A

X

B

Single RIB Route Server

C

D

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

# Single-RIB BGP policy problem

**IXP Fabric**

Y ↔ A

X

B

**Single RIB Route Server**

C

D

- ☑ ABCD route server clients
- ☑ X reaches IXP via transit A and B
- ☑ AYX longer than BX

# Single-RIB BGP policy problem



**IXP Fabric**

Y

A

X

B

**Single RIB Route Server**

C

D

☑ RS calculates best path as BX

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

# Single-RIB BGP policy problem

**IXP Fabric**

**Single RIB Route Server**

Y

A

X

B

C

D

☑ RS calculates best path as BX

☑ B does not peer with D (tag `0:D`)

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

# Single-RIB BGP policy problem

**IXP Fabric**

Y

A

X

B

**Single RIB Route Server**

C

D

☑ RS calculates best path as BX

☑ B does not peer with D (tag `0:D`)

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

☑ D should see XYA path

# Single-RIB BGP policy problem

**IXP Fabric**

**Y**

**A**

**X**

**B**

**Single RIB Route Server**

**C**

**D**

☑ RS calculates best path as BX

☑ B does not peer with D (tag `0:D`)

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

☑ D should see XYA path

☑ RS RIB has only one best path - BX

# Single-RIB BGP policy problem

**IXP Fabric**

Y ⟷ A

Y ⟷ X

A ⟷ Single RIB Route Server

X ⟷ B

B ⟷ Single RIB Route Server

Single RIB Route Server ⟷ C

Single RIB Route Server ⟷ D

☑ RS calculates best path as BX

☑ B does not peer with D (tag `0:D`)

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

☑ D should see XYA path

☑ RS RIB has only one best path - BX

☑ D does not see X

# How Per-Client Loc-RIBs Work



export
import

Y
A
X
B

Multi-RIB Route Server

Loc-RIB A
Loc-RIB B
Loc-RIB C
Loc-RIB D

C
D

# How Per-Client Loc-RIBs Work

export
import

Y

A

X

Multi-RIB Route Server

Loc-RIB A

Loc-RIB B

Loc-RIB C

Loc-RIB D

C

B

D

☑ ABCD route server clients

export
import

Multi-RIB Route Server

Loc-RIB A

Loc-RIB B

Loc-RIB C

Loc-RIB D

Y

A

X

C

B

D

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

# How Per-Client Loc-RIBs Work

export
import

Multi-RIB Route Server

Loc-RIB A
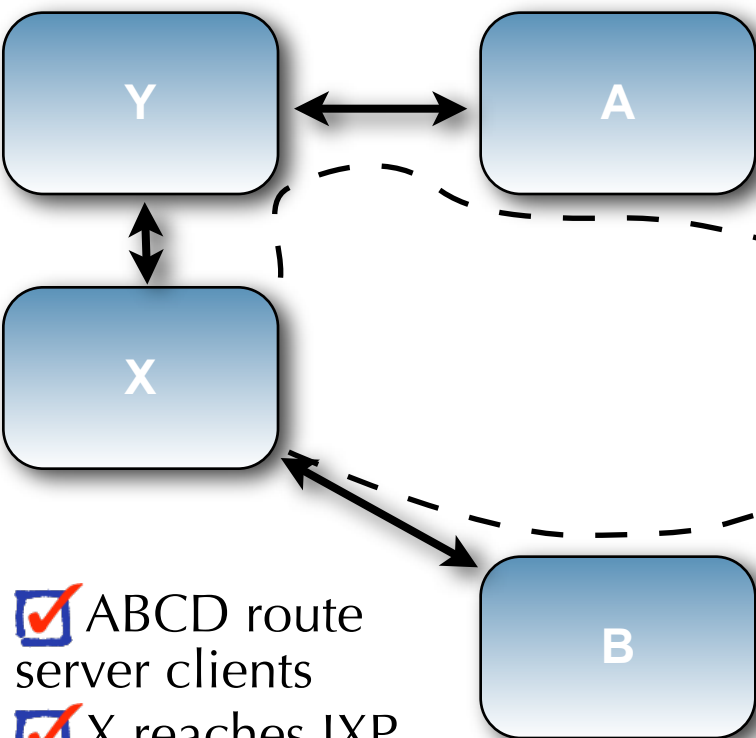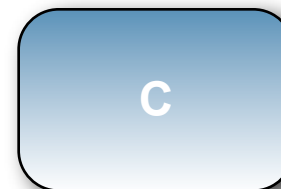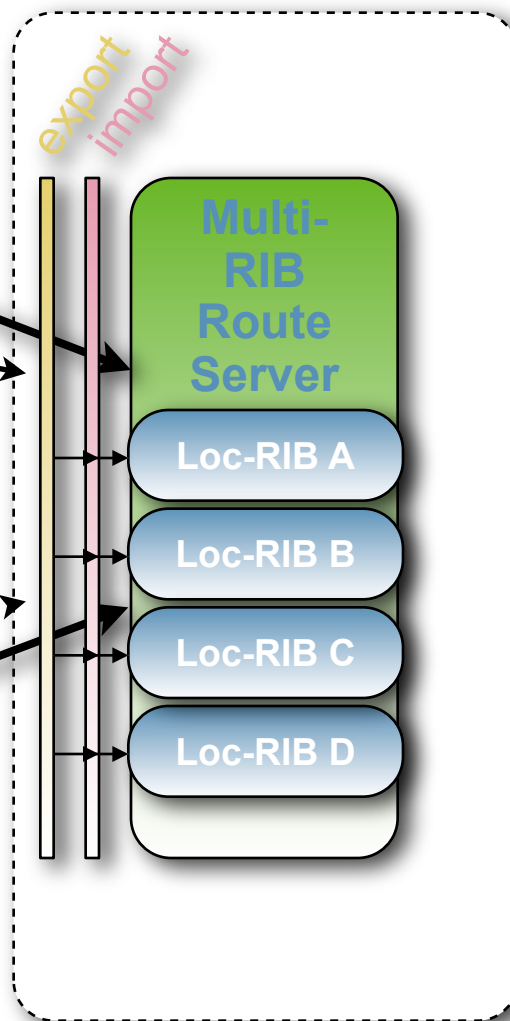
Loc-RIB B

Loc-RIB C

Loc-RIB D

Y

A

X

B

C

D

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

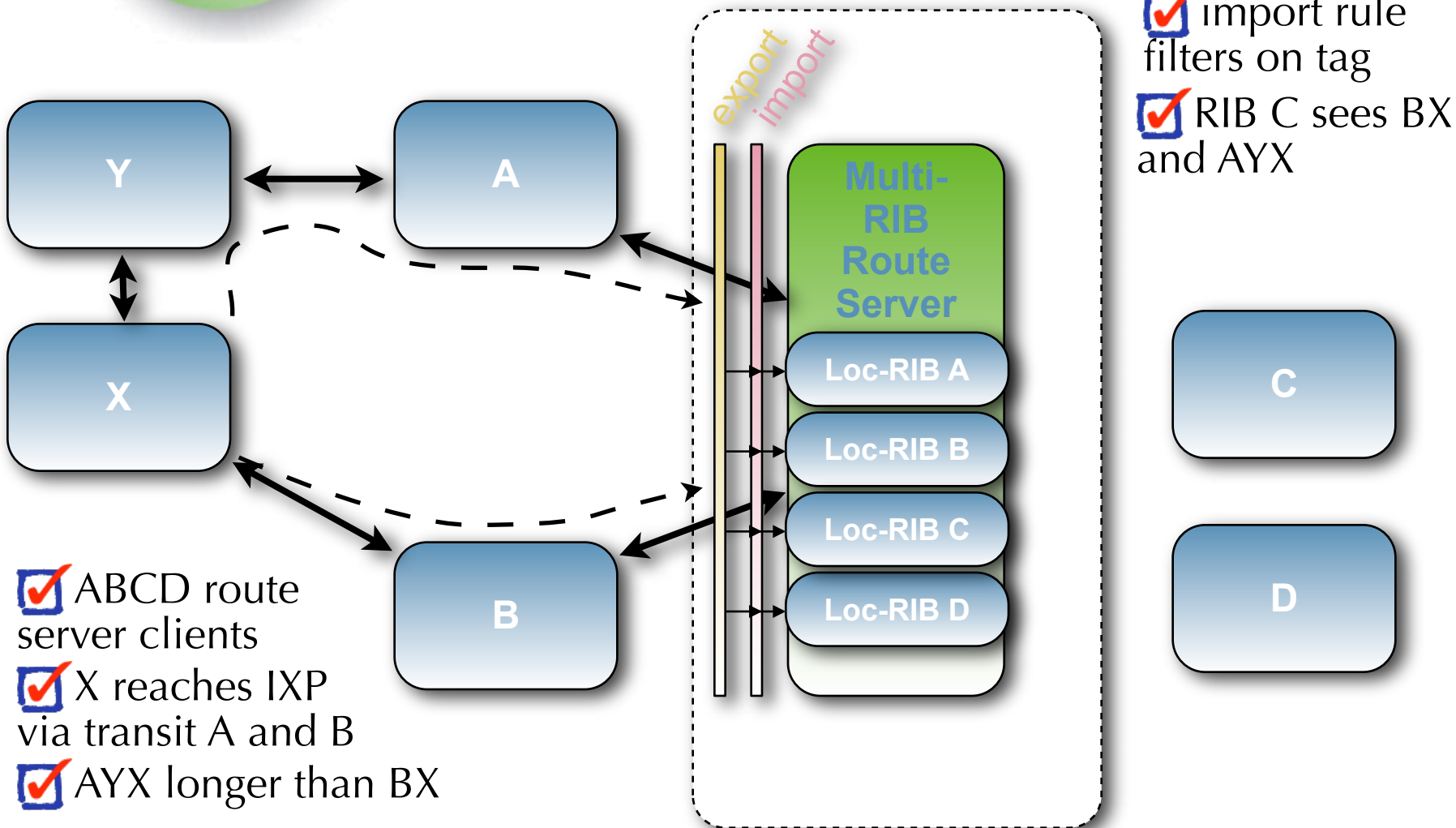# How Per-Client Loc-RIBs Work



☑ import rule filters on tag

export

import

Y

A

X

Multi-RIB Route Server

Loc-RIB A

Loc-RIB B

Loc-RIB C

Loc-RIB D

C

B

D

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

# How Per-Client Loc-RIBs Work

**export** **import**

**Multi-RIB Route Server**

Loc-RIB A

Loc-RIB B

Loc-RIB C

Loc-RIB D

Y

A

X

B

C

D

☑ import rule filters on tag

☑ RIB C sees BX and AYX

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

# How Per-Client Loc-RIBs Work

export  import

Y ⟷ A

Y ⟷ X

X

B

**Multi-RIB Route Server**

Loc-RIB A

Loc-RIB B

Loc-RIB C

Loc-RIB D

☑ import rule filters on tag

☑ RIB C sees BX and AYX

☑ RIB D only sees AYX

C

D

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

# How Per-Client Loc-RIBs Work



import rule filters on tag

RIB C sees BX and AYX

RIB D only sees AYX

**Multi-RIB Route Server**

Loc-RIB A

Loc-RIB B

Loc-RIB C

Loc-RIB D

export

import

BX

Y

A

X

C

B

D

ABCD route server clients

X reaches IXP via transit A and B

AYX longer than BX

RIB C selects BX

# How Per-Client Loc-RIBs Work

export   import

**Multi-RIB Route Server**

Loc-RIB A

Loc-RIB B

Loc-RIB C

Loc-RIB D

Y

A

X

C

B

D

BX

AYX

☑ import rule filters on tag

☑ RIB C sees BX and AYX

☑ RIB D only sees AYX

☑ ABCD route server clients

☑ X reaches IXP via transit A and B

☑ AYX longer than BX

☑ RIB C selects BX

☑ RIB D selects AYX

Ceiling Cat is Watching you Propagate

- Multiple Loc-RIBs mean:
  - Memory, CPU consumption go from $O(M)$ to $O(N \times M)$
    - $N$ = number of clients
    - $M$ = total number of prefixes

- Multiple Loc-RIBs mean:
  - Memory, CPU consumption go from O(M) to O(N x M)
    - N = number of clients
    - M = total number of prefixes
  - Update processing resources required are:
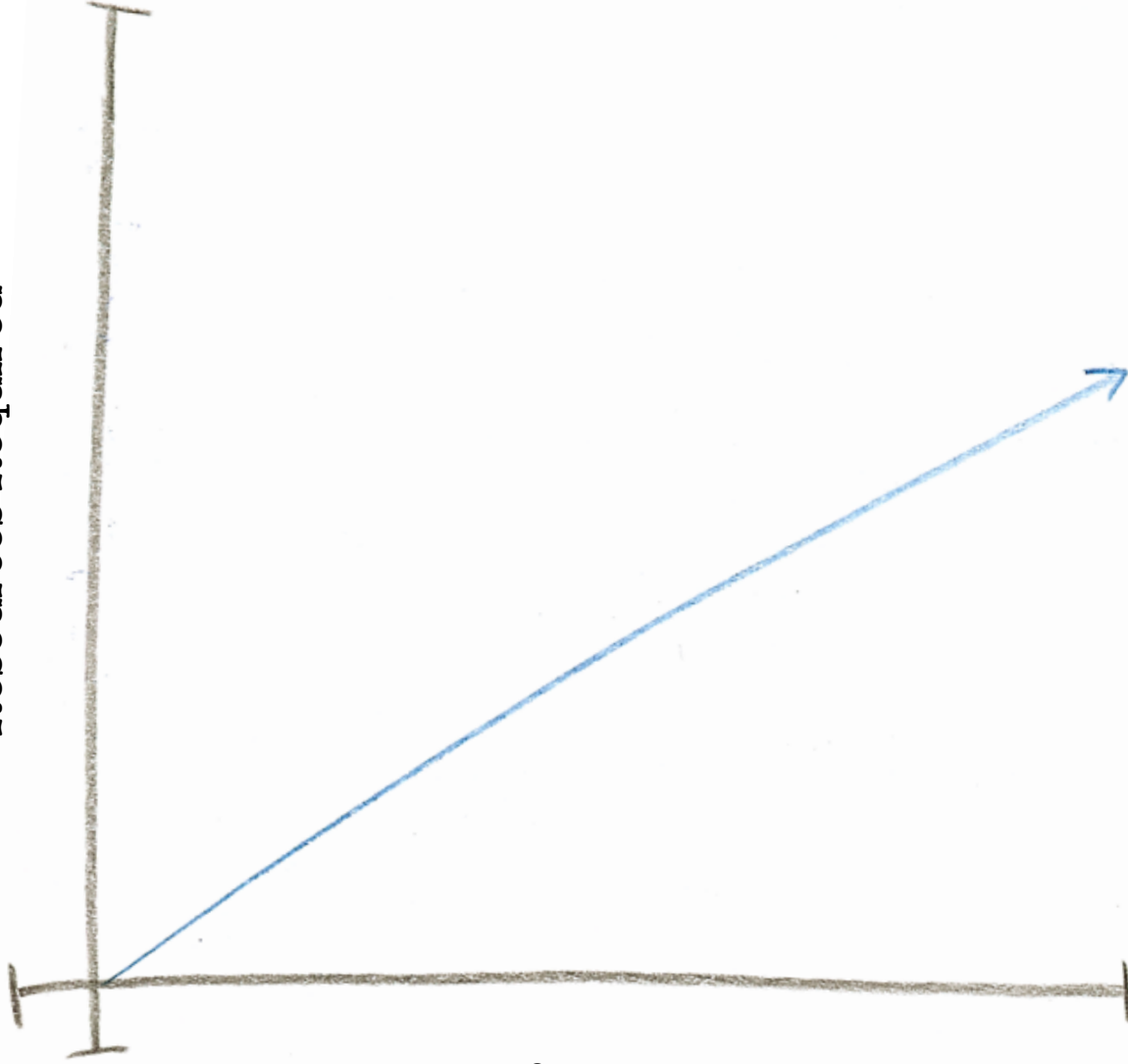
$$\sum_{1}^{N}(P(n) \bullet (N-1))$$

  - Where
    - P(n) = the number of prefixes from peer n
    - N = number of peers connected to system
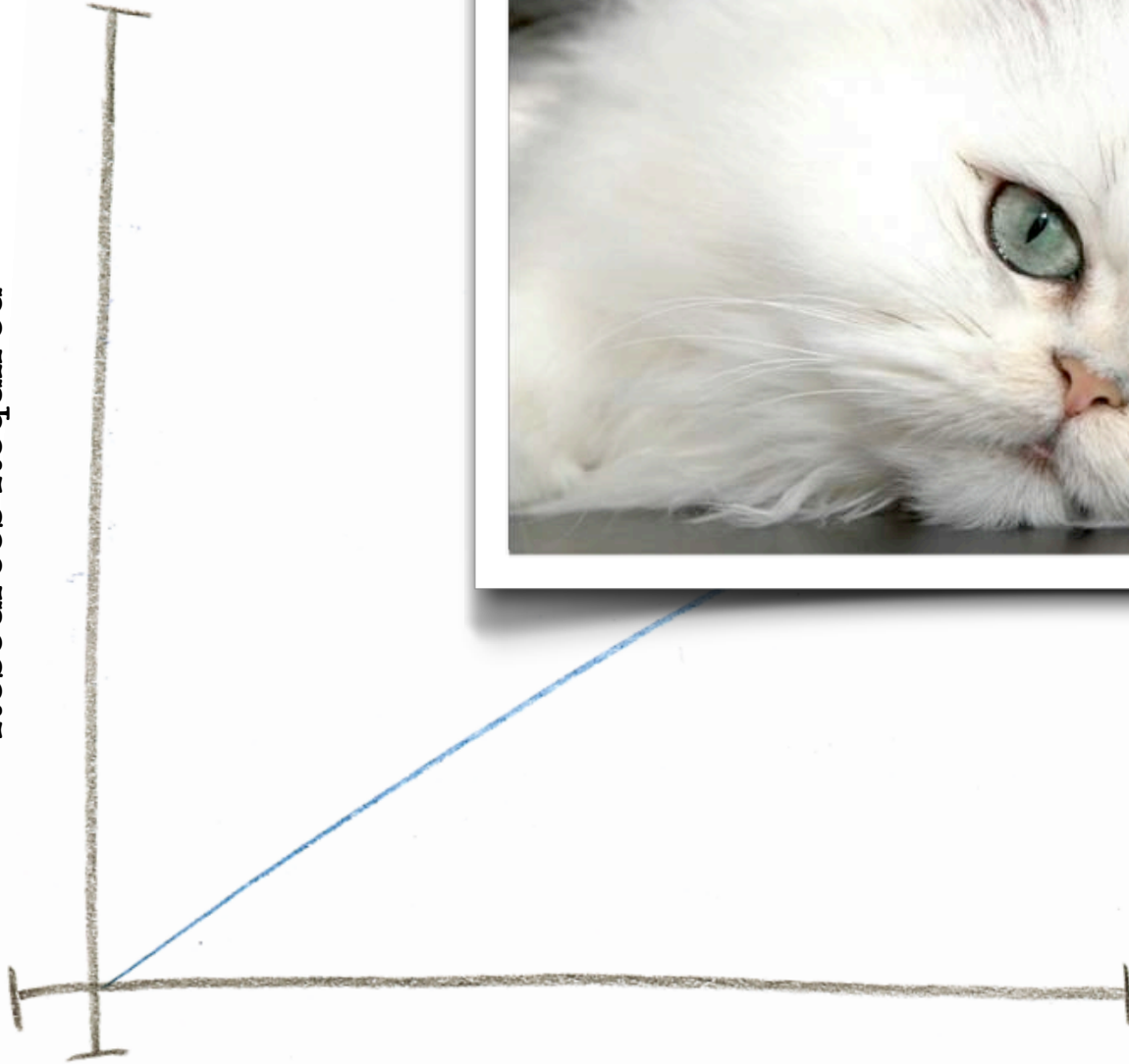  - This scales as $P_{average} * N^2$
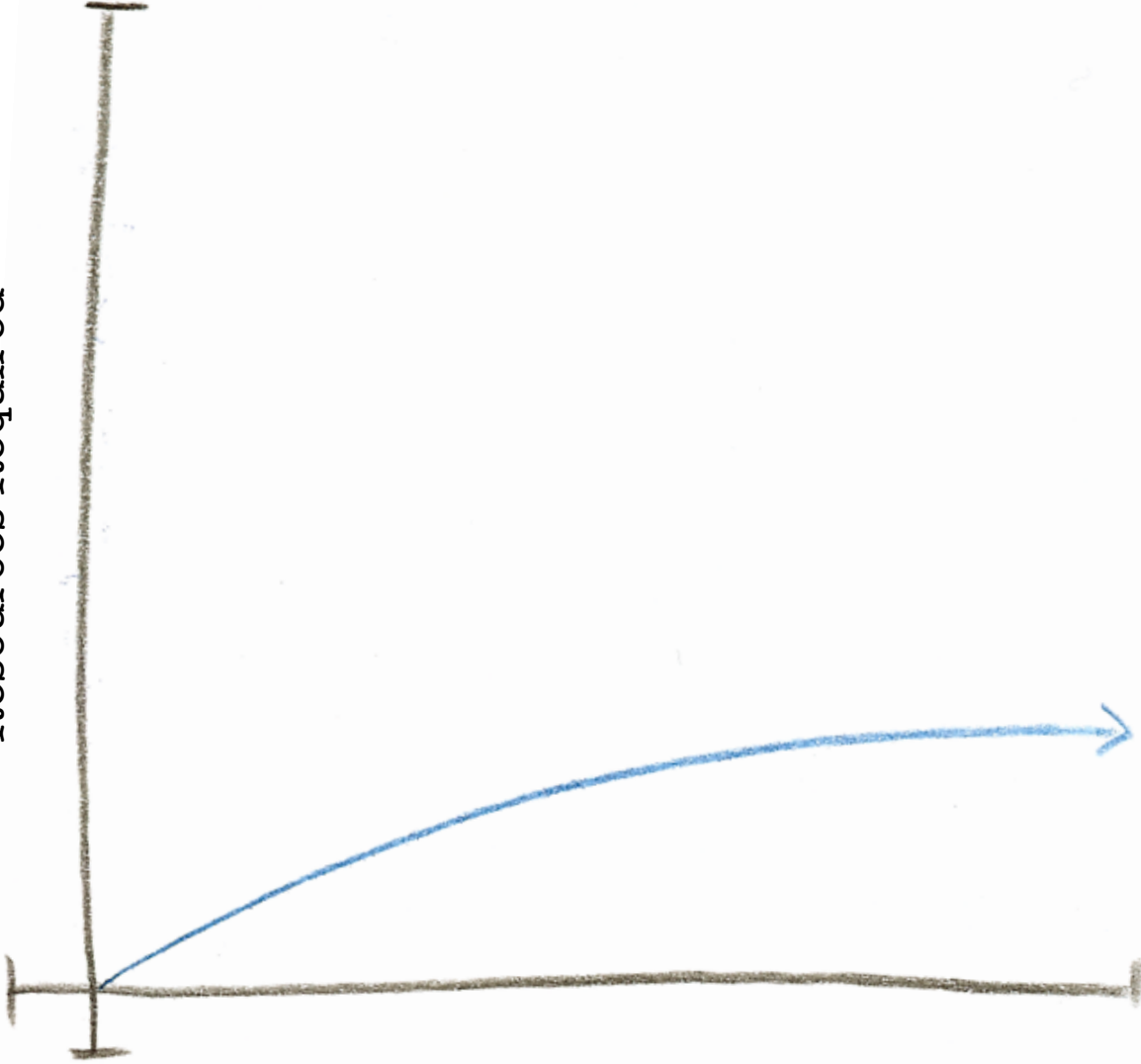
Cute Kitteh!



Resources Required

Peers & Prefixes

Resources Required
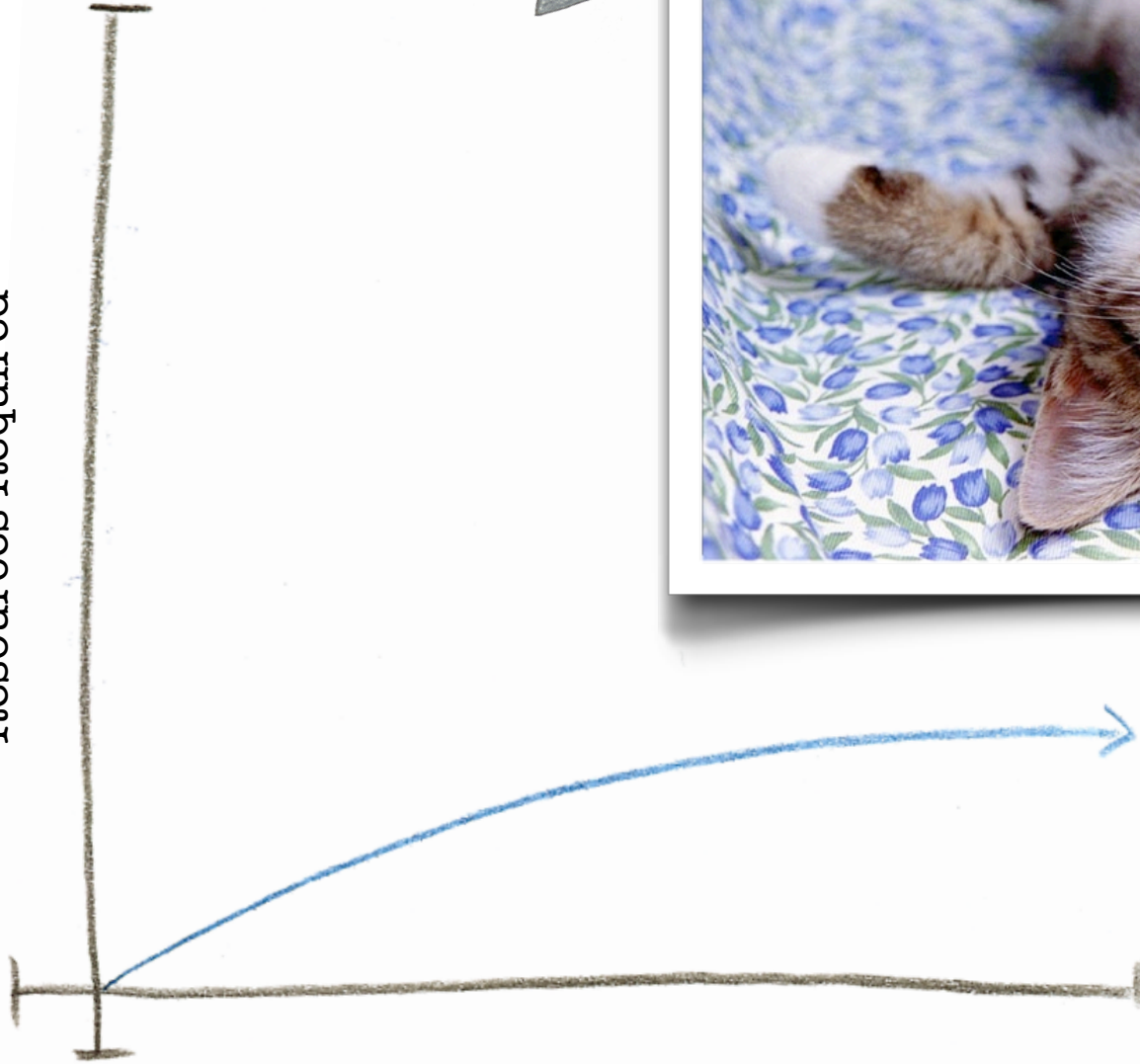
Peers & Prefixes

# SRSLY Cute Kitteh!



Resources Required

Peers & Prefixes

Resources Required vs. Peers & Prefixes

Evil Kitteh!

Resources Required
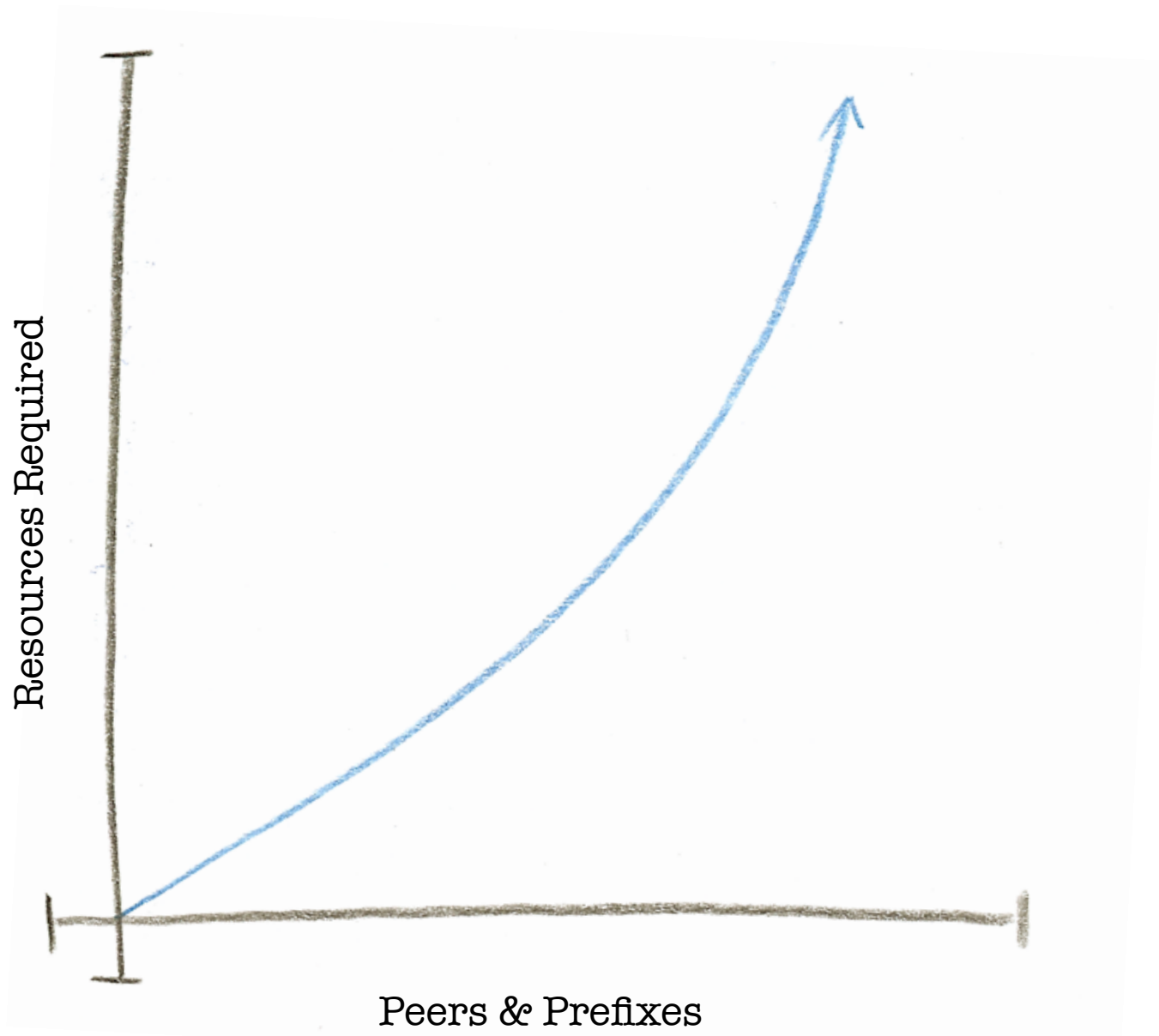
Peers & Prefixes

Resources Required

Peers & Prefixes

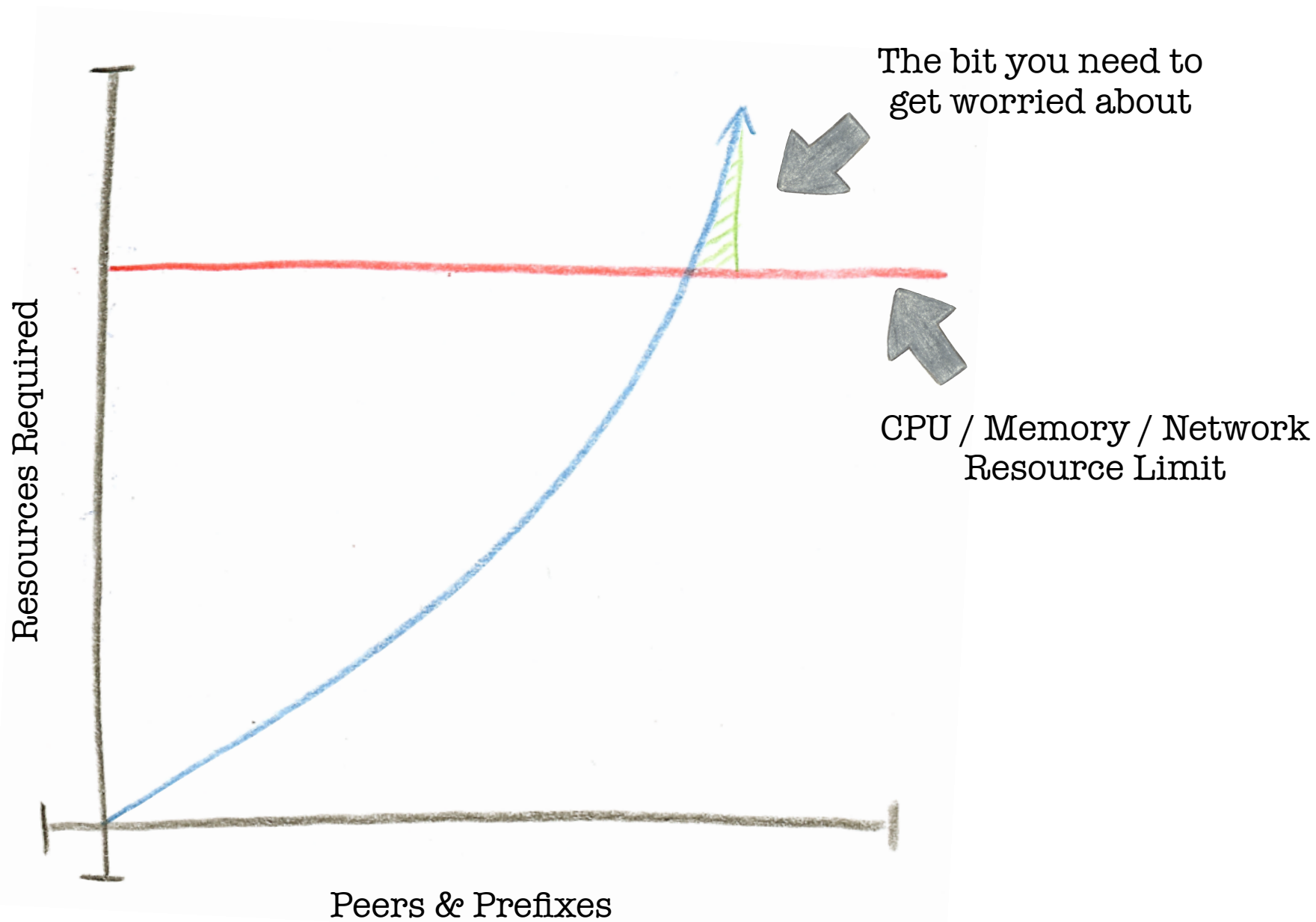Resources Required

Peers & Prefixes

CPU / Memory / Network
Resource Limit

The bit you need to get worried about

CPU / Memory / Network Resource Limit

Resources Required

Peers & Prefixes

- A single BGP prefix update
  - might take 10 - 30 bytes on network to send to peer
  - might take 10 - 30μS to process update
- Disclaimer
  - this ignores attributes, path length, cpu speed, and a pile of other highly relevant parameters

- A single BGP prefix update
  - might take 10 - 30 bytes on network to send to peer
  - might take 10 - 30μS to process update
- Disclaimer
  - this ignores attributes, path length, cpu speed, and a pile of other highly relevant parameters
- Ok, it's hand-waving

# Facts and Figures

- 200 clients
- Average 100 prefixes each

# Facts and Figures



- 200 clients
- Average 100 prefixes each

**Multi-RIB Route Server**
- RIB A
- RIB B
- RIB C
- ⋮
- RIB 100

19500 prefixes → Client A
50 prefixes ←

19500 prefixes → Client B
50 prefixes ←

19500 prefixes → Client C
50 prefixes ←

⋮

19500 prefixes
50 prefixes → Client 100

**Multi-RIB Route Server**

RIB A
RIB B
RIB C
⋮
RIB 100

19500 prefixes
50 prefixes

19500 prefixes
50 prefixes

19500 prefixes
50 prefixes

19500 prefixes
50 prefixes

*Client A*

*Client B*

*Client C*

⋮

*Client 100*

- 200 clients
- Average 100 prefixes each

- RS to Client updates:
  - 19500*200 = 3,900,000
- Client to RS updates:
  - 100 * 200 = 20000 updates
- Total BGP updates: 4,000,000

- 40-120M of network traffic
- 40-120s CPU time

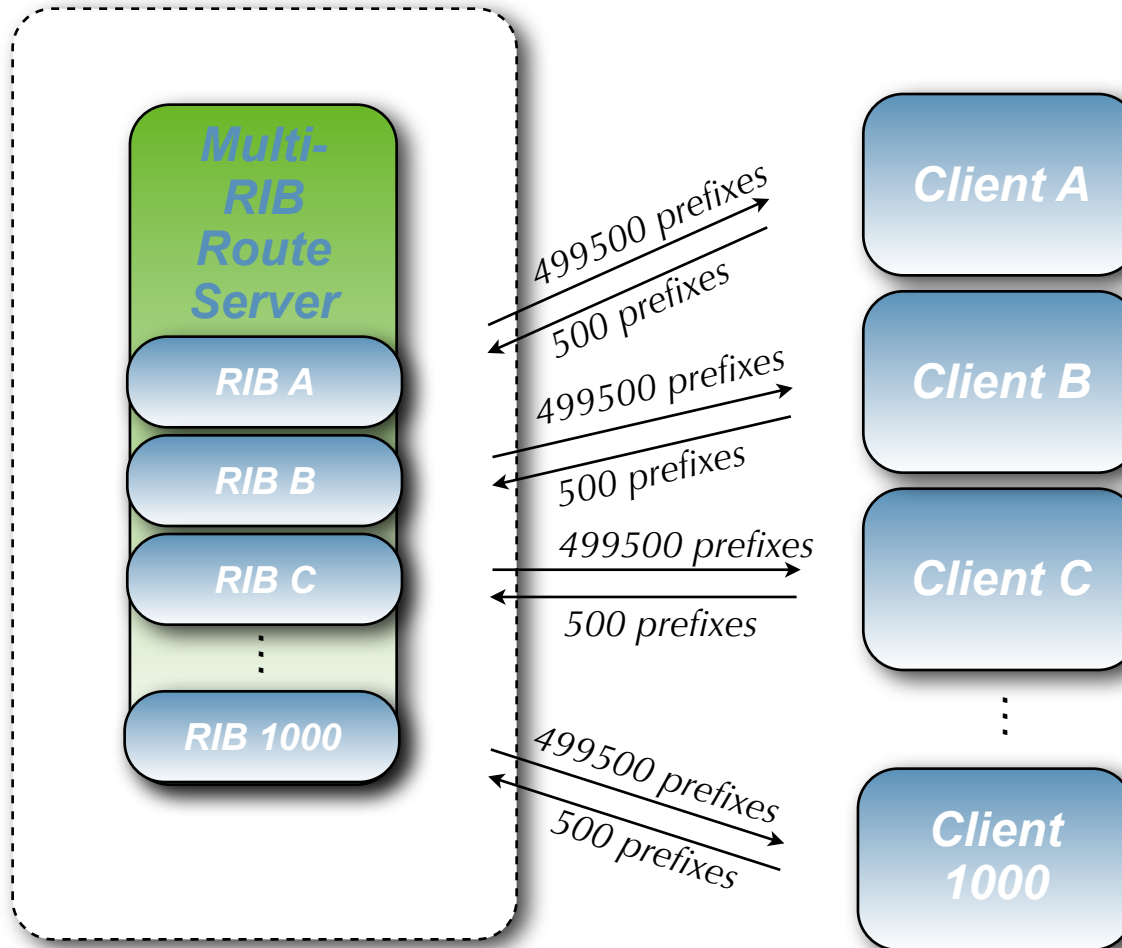## Facts and Figures

- 1000 clients
- Average 500 prefixes each

# Facts and Figures

- 1000 clients
- Average 500 prefixes each

**Multi-RIB Route Server**

- RIB A
- RIB B
- RIB C
- ⋮
- RIB 1000

Client A — 499500 prefixes → / 500 prefixes ←

Client B — 499500 prefixes → / 500 prefixes ←

Client C — 499500 prefixes → / 500 prefixes ←

⋮

Client 1000 — 499500 prefixes → / 500 prefixes →

# Facts and Figures

- 1000 clients
- Average 500 prefixes each

- RS to Client updates:
  - 499000*1000 = 499.5m
- Client to RS updates:
  - 1000 * 500 = 500k updates
- Total BGP updates: 500m

- 5-15G of network traffic
- 85 - 250m CPU time

**Multi-RIB Route Server**

RIB A

RIB B

RIB C

RIB 1000

Client A

Client B

Client C

Client 1000

499500 prefixes
500 prefixes
499500 prefixes
500 prefixes
499500 prefixes
500 prefixes
499500 prefixes
500 prefixes

- Super-linear scaling causes inherent breakage
  - Moving away from one Loc-RIB per client model is critical
  - Right now, this isn't the primary cause of IXP breakage
- Three primary models to escape this limitation
  - Collapse multiple Loc-RIBs in memory into single gargantuan Loc-RIB
    - less memory, less CPU
    - "You can run but you can't hide"
  - Use prior knowledge
    - "Web based peering"
    - Disable unique Loc-RIB on per client basis
  - BGP `ADD_PATH` Capability
    - Published as ID: `draft-walton-bgp-add-paths`
    - Moves BGP Best Path Selection to client, so filtering can be performed without selecting