# The Evolution of Traffic Routing in a Streaming World

Abhishek Srikanth

SOFTWARE ENGINEER,
FACEBOOK

- Purdue University

- R&D Intern @ Bloomberg

- Program Manager Intern @ Microsoft

- Software Engineer @ Facebook
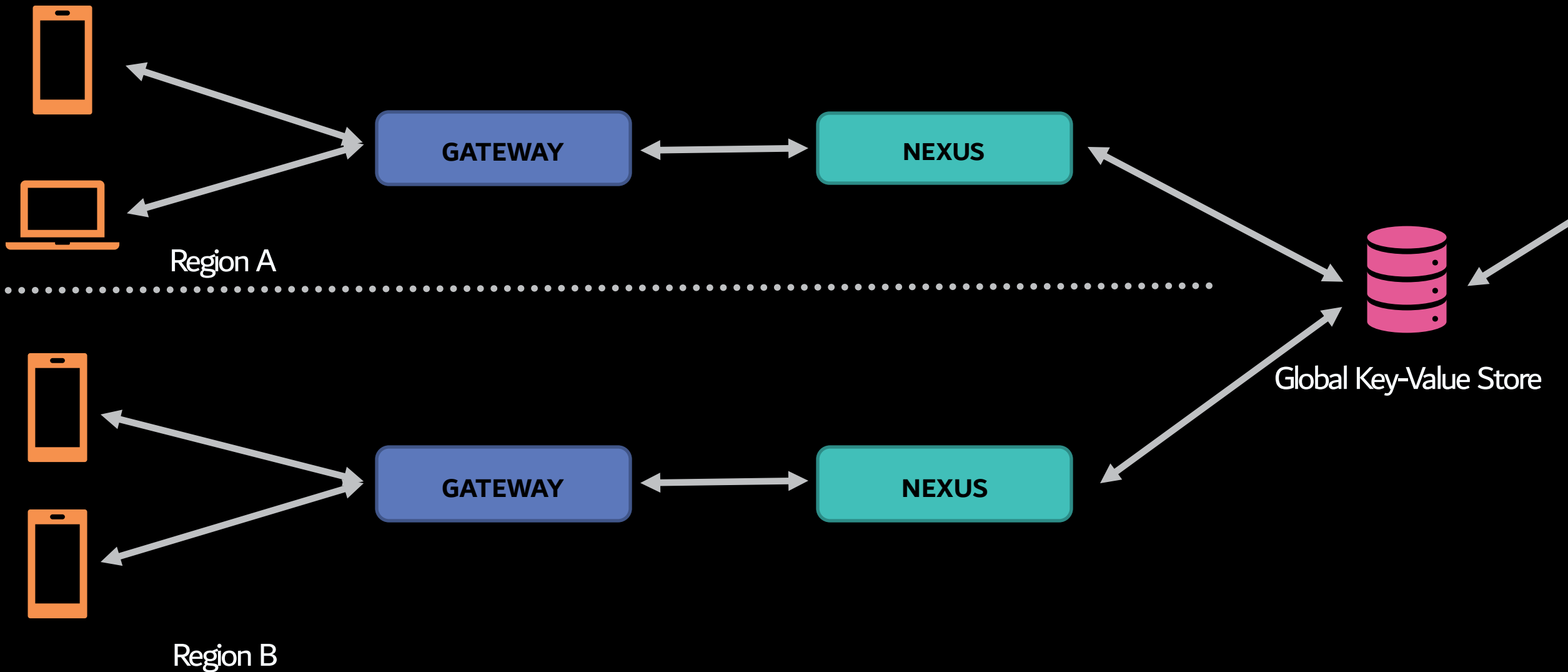
A little about me

# The Evolution

Collocated Architecture

Off-box Architecture

Full-Mesh Architecture

5000 ft view

GATEWAY — NEXUS

GATEWAY — NEXUS

Global Key-Value Store

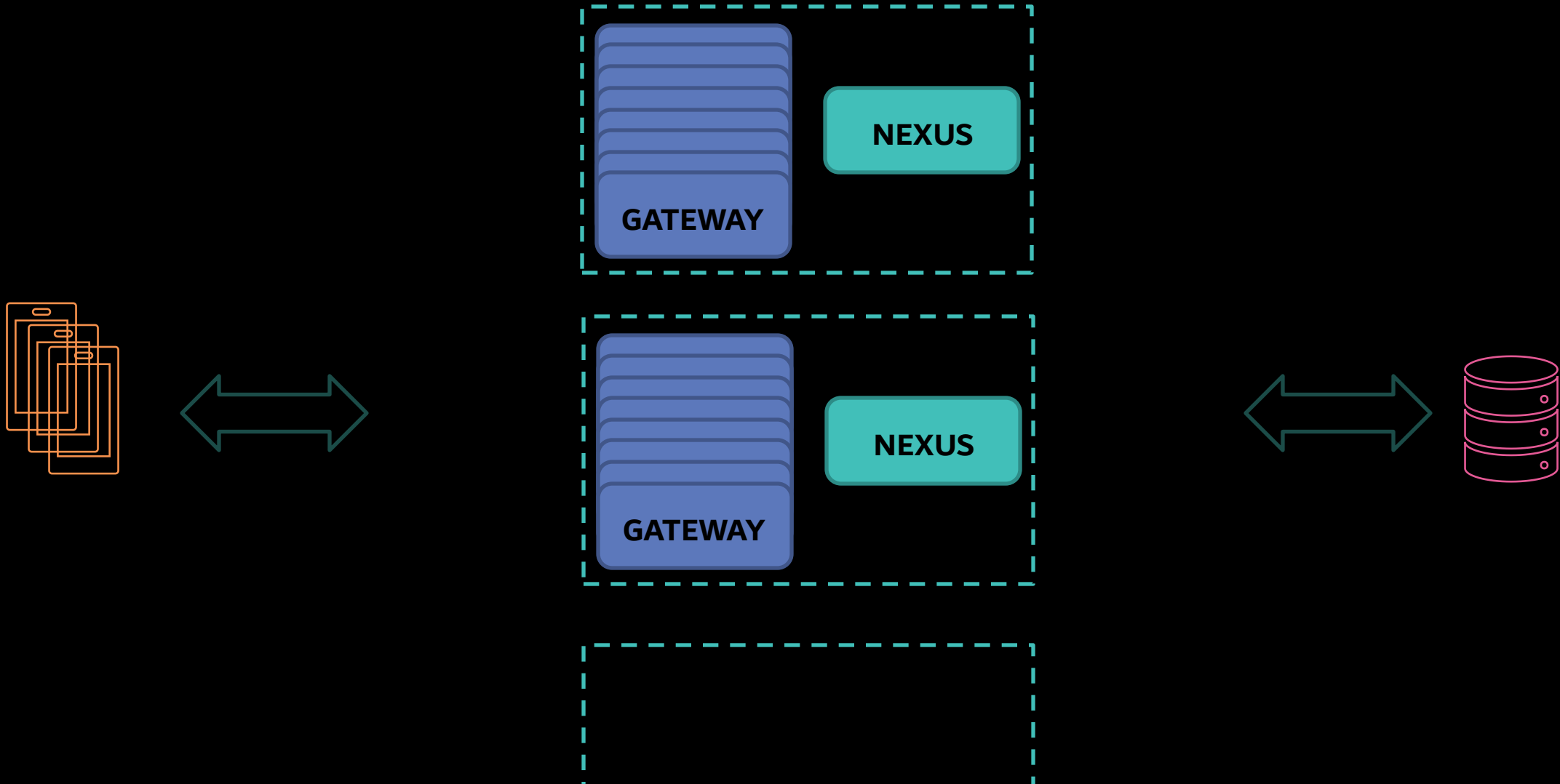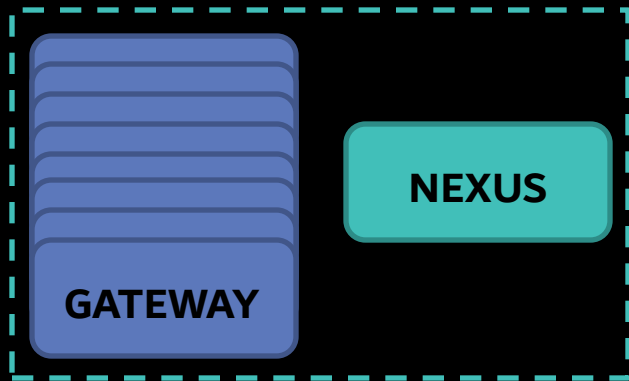# Collocated Architecture

# Collocated Architecture

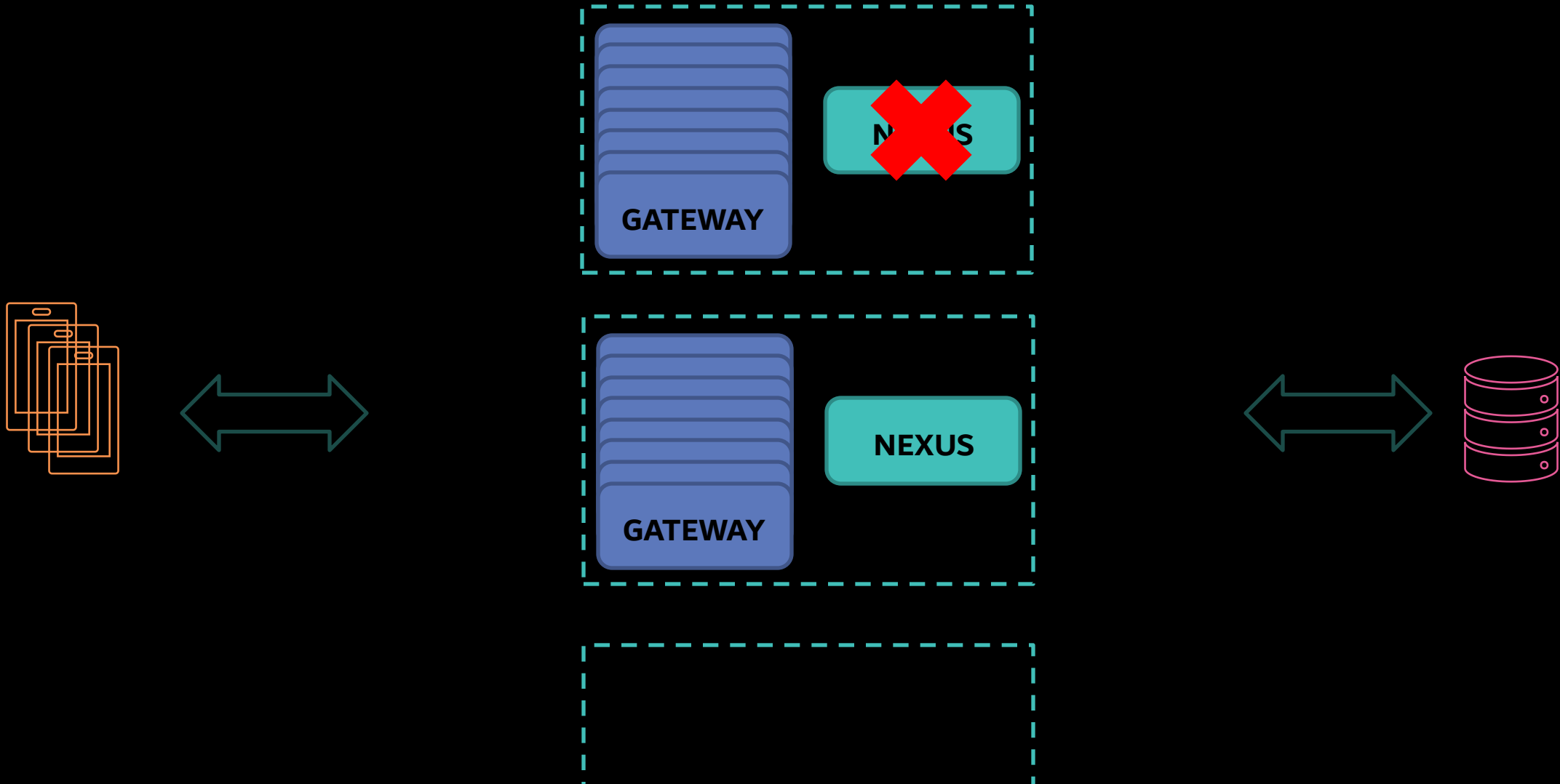# Collocated Architecture

1. Shared Resources

2. Independent Deployments

3. Low Fault Tolerance

# Collocated Architecture

# Collocated Architecture

# Off-box Architecture

# Off-box Architecture

# Off-box Architecture

1. Customer Traffic Isolation

2. High Fault Tolerance

3. Dynamic Load Balancing

# Off-box Architecture

Advantages : Customer Traffic Isolation



Isolation for customer A (E.g. live videos)

Shared isolation for all other customers

# Off-box Architecture

Advantages : High Fault Tolerance

# Off-box Architecture

Advantages : High Fault Tolerance

# Off-box Architecture

Advantages : High Fault Tolerance

# Off-box Architecture

Advantages : Dynamic Load Balancing

# Off-box Architecture

Advantages : Dynamic Load Balancing

# Off-box Architecture

Advantages : Dynamic Load Balancing

# Off-box Architecture

Advantages : Dynamic Load Balancing

# Off-box Architecture

Advantages : Dynamic Load Balancing

# Off-box Architecture

## Challenges: Socket Connections



In a Full-Mesh connection:

Total Number of Connections: G * N

Number of Connections per Nexus: G

# Off-box Architecture

## Challenges: Socket Connections

**GATEWAY**

**GATEWAY**

**GATEWAY**

**GATEWAY**

**GATEWAY**

**GATEWAY**

**NEXUS**

**NEXUS**

**NEXUS**

**NEXUS**

**NEXUS**

**NEXUS**

In a Full-Mesh connection:

Total Number of Connections: $G * N$

Number of Connections per Nexus: $G$

If each Gateway Connected to only 'K' Nexuses:

Total Number of Connections: $G * K$

Number of Connections per Nexus: $G * K / N$

For our system, K = 3 provided desirable results.

# Off-box Architecture

# Off-box Architecture

## Challenges: Which Nexus to talk to?

# Off-box Architecture

Challenges: Which Nexus to talk to? A Naïve Solution

# Off-box Architecture

Challenges: Which Nexus to talk to? A Naïve Solution

# Off-box Architecture

GATEWAY    GATEWAY    GATEWAY

NEXUS    NEXUS    NEXUS    NEXUS    NEXUS    NEXUS    NEXUS

Unavailable
(Due to deployment)

# Off-box Architecture

Available Again

# Off-box Architecture

Challenges: Which Nexus to talk to?

Conditions:

- Independent Gateway Decisions

- Even Distribution of Connections

- Minimal Disruptions

# Rendezvous Hash

## What is it?

Rendezvous hashing is an algorithm that allows clients to achieve distributed agreement on a set of 'k' options out of a possible set of 'n' options.

# Rendezvous Hash

## What is it?

Rendezvous hashing is an algorithm that allows clients to achieve distributed agreement on a set of 'k' options out of a possible set of 'n' options.

# Rendezvous Hash

## How does it work?

On Gateway $g_i$:

For $n_j$ in N:

   weight_of_$n_j$ = HASH( $g_i$ , $n_j$ )

Connect to nexuses with the top 'k' weights values

# Rendezvous Hash

How does it work?

On Gateway $g_i$:

For $n_j$ in N:

    weight_of_$n_j$ = HASH( $g_i$ , $n_j$ )

Connect to nexuses with the top 'k' weights values



$g_1$        $g_2$        $g_3$        $g_G$

GATEWAY    GATEWAY    GATEWAY        GATEWAY

?    ?    ?

'k' = 3

NEXUS    NEXUS    NEXUS    NEXUS       NEXUS

78   $n_1$     13   $n_2$     44   $n_3$     9   $n_4$       15   $n_N$

# Rendezvous Hash

How does it work?

On Gateway $g_i$:
For $n_j$ in N:
   weight_of_$n_j$ = HASH( $g_i$ , $n_j$ )
Connect to nexuses with the top 'k' weights values

# Rendezvous Hash

How does it work?

On Gateway $g_i$:
For $n_j$ in N:
    weight_of_$n_j$ = HASH( $g_i$ , $n_j$ )
Connect to nexuses with the top 'k' weights values

# Rendezvous Hash

## A Note on the Hash and Distribution

On Gateway g:
For $n_j$ in N:
    weight_of_$n_j$ = HASH( $g_i$ , $n_j$ )
Connect to nexuses with the top 'k' weights values

# Rendezvous Hash

Minimal Disruption



|  | NEXUS | NEXUS | NEXUS | NEXUS | NEXUS | NEXUS |
|---|---|---|---|---|---|---|
| $g_1$ | 88 | 12 | 37 | 78 | 45 | 58 |
| $g_2$ | 99 | 98 | 3 | 9 | 21 | 6 |
| $g_3$ | 6 | 23 | 52 | 26 | 49 | 45 |

# Rendezvous Hash

Minimal Disruption



|       |    |    |    |    |    |    |    |
|-------|----|----|----|----|----|----|----|
| g₁    | 88 | 12 | 37 | 78 | 45 | 58 | 3  |
| g₂    | 99 | 98 | 3  | 9  | 21 | 6  | 18 |
| g₃    | 6  | 23 | 52 | 26 | 49 | 45 | 73 |

# Rendezvous Hash

Minimal Disruption



| | NEXUS | NEXUS | NEXUS | NEXUS | NEXUS | NEXUS | NEXUS |
|-----|-----|-----|-----|-----|-----|-----|-----|
| $g_1$ | 88 | 12 | 37 | 78 | 45 | 58 | 3 |
| $g_2$ | 99 | 98 | 3 | 9 | 21 | 6 | 18 |
| $g_3$ | 6 | 23 | 52 | 26 | 49 | 45 | 73 |

# Rendezvous Hash

Minimal Disruption



|       |    |    |    |    |     |    |    |
|-------|----|----|----|----|-----|----|----|
| $g_1$ | 88 | 12 | 37 | 78 | 45  | 58 | 3  |
| $g_2$ | 99 | 98 | 3  | 9  | 21  | 6  | 18 |
| $g_3$ | 6  | 23 | 52 | 26 | 49  | 45 | 73 |

# Rendezvous Hash

Minimal Disruption



|  | | | | | | |
|---|---|---|---|---|---|---|
| $g_1$ | 88 | 12 | 37 | 78 | 45 | 58 | 3 |
| $g_2$ | 99 | 98 | 3 | 9 | 21 | 6 | 18 |
| $g_3$ | 6 | 23 | 52 | 26 | 49 | 45 | 73 |

# Off-box Architecture

# Off-box Architecture

Traffic Routing

# Off-box Architecture

Traffic Routing

# Off-box Architecture

Traffic Routing

# Off-box Architecture

Traffic Routing

# Off-box Architecture

# Off-box Architecture

# Off-box Architecture

Traffic Routing: Lower of random two hosts



GATEWAY

After time 't'

GATEWAY

16.7%

66.6%

16.7%

NEXUS — Streams: 100    NEXUS — Streams: 5    NEXUS — Streams: 100

NEXUS — Streams: 134    NEXUS — Streams: 137    NEXUS — Streams: 134

# Off-box Architecture

1. Customer Traffic Isolation

2. High Fault Tolerance

3. Dynamic Load Balancing

# Full-Mesh Architecture

5000 ft view

GATEWAY ⟷ WARPGATE ⟷ BLADERUNNER

GATEWAY ⟷ WARPGATE ⟷ BLADERUNNER

Global Key-Value Store

# Full-Mesh Architecture

# Full-Mesh Architecture

Sticky Routing

# Full-Mesh Architecture

Sticky Routing

# Full-Mesh Architecture

## Sticky Routing

# Full-Mesh Architecture

15-20% reduction of requests to the Global Key-Value Store

5-10% Memory and CPU Savings on Bladerunner

# Full-Mesh Architecture

Sticky Routing

How?

# Full-Mesh Architecture

Sticky Routing : Rendezvous Hashing

# Full-Mesh Architecture

On any WarpGate:
Say we are routing feed $f_i$

# Full-Mesh Architecture

## Sticky Routing : Rendezvous Hashing

$f_1$

On any WarpGate:
Say we are routing feed $f_i$
For $br_j$ in Bladerunner:
    weight_of_$br_j$ = HASH( $f_i$ , $br_j$ )

Finally, route to Bladerunner with highest weight.



WARPGATE

WARPGATE

WARPGATE

WARPGATE

WARPGATE

WARPGATE

BLADERUNNER    27

BLADERUNNER    73

BLADERUNNER    64

BLADERUNNER    81

BLADERUNNER    7

# Full-Mesh Architecture

Sticky Routing : Rendezvous Hashing

$f_1$

On any WarpGate:

Say we are routing feed $f_i$

For $br_j$ in Bladerunner:

  weight_of_$br_j$ = HASH( $f_i$ , $br_j$ )

Finally, route to Bladerunner with highest weight.

WARPGATE

WARPGATE

WARPGATE

WARPGATE

WARPGATE

WARPGATE

BLADERUNNER    27

BLADERUNNER    73

BLADERUNNER    64

BLADERUNNER    81

BLADERUNNER    7

# Full-Mesh Architecture

## Sticky Routing

On any WarpGate:

Say we are routing feed $f_i$

For $br_j$ in Bladerunner:

$weight\_of\_br_j = HASH(\ f_i\ ,\ br_j\ )$

Finally, route to Bladerunner with highest weight.

| | $f_1$ | $f_2$ |
|---|---|---|
| BLADERUNNER | 27 | 17 |
| BLADERUNNER | 73 | 54 |
| BLADERUNNER | 64 | 32 |
| BLADERUNNER | 81 | 23 |
| BLADERUNNER | 7 | 2 |

# Full-Mesh Architecture

## Sticky Routing : Minimal Disruption

# Full-Mesh Architecture

## Sticky Routing : Minimal Disruption

# Full-Mesh Architecture

## Sticky Routing

# Full-Mesh Architecture

WARPGATE
WARPGATE
WARPGATE
WARPGATE
WARPGATE
WARPGATE

BLADERUNNER
BLADERUNNER
BLADERUNNER
BLADERUNNER

1. Sticky Routing

2. Virtual Isolation

# Looking back

Customer traffic Isolation in the off-box architecture



GATEWAY    GATEWAY    GATEWAY    GATEWAY    GATEWAY

NEXUS    NEXUS    NEXUS    NEXUS    NEXUS    NEXUS

Isolation for customer A (E.g. live videos)

Shared isolation for all other customers

# Full-Mesh Architecture

| Bladerunner 1 | Lorem | Ipsum | Dolor | Sit | Amet |
|---|---|---|---|---|---|
| Bladerunner 2 | Lorem | Ipsum | Dolor | Sit | Amet |
| Bladerunner 3 | Lorem | Ipsum | Dolor | Sit | Amet |
| Bladerunner 4 | Lorem | Ipsum | Dolor | Sit | Amet |
| Bladerunner 5 | Lorem | Ipsum | Dolor | Sit | Amet |
| Bladerunner 6 | Lorem | Ipsum | Dolor | Sit | Amet |
| Bladerunner 7 | Lorem | Ipsum | Dolor | Sit | Amet |
| ... | | | | | |
| Bladerunner 'b' | Lorem | Ipsum | Dolor | Sit | Amet |

# Full-Mesh Architecture

Virtual Isolation

| | |
|---|---|
| **Bladerunner 1** | Lorem |
| **Bladerunner 2** | Ipsum |
| **Bladerunner 3** | Dolor |
| **Bladerunner 4** | Ipsum |
| **Bladerunner 5** | Sit |
| **Bladerunner 6** | Lorem |
| **Bladerunner 7** | Amet |
| **...** | |
| **Bladerunner 'b'** | Dolor |

# Full-Mesh Architecture

Virtual Pools for Isolation

Noisy Neighbour

Resource Attribution

Maintenance Costs

# Full-Mesh Architecture

Virtual Pools for Isolation

1. WarpGate Visibility

# Full-Mesh Architecture

Virtual Pools for Isolation: WarpGate Visibility

WARPGATE   WARPGATE   WARPGATE   WARPGATE   WARPGATE   WARPGATE

HOST LB FACTOR: 100000,  HEAT STATUS: NORMAL,
CUSTOMER'S HANDLER: {

  IPSUM: {    LB FACTOR: 35000,    HEAT STATUS: LOW_UTIL },
  AMET:  {    LB FACTOR: 65000,    HEAT STATUS: WARM }, }

**Bladerunner 4**   Ipsum   Amet

# Full-Mesh Architecture

Virtual Pools for Isolation: WarpGate Visibility

# Full-Mesh Architecture

1. WarpGate Visibility

2. Virtual Pool placement

# Full-Mesh Architecture

Virtual Pools for Isolation: Virtual Pool Placement

On each WarpGate:

Say we are to scale up vpool $_i$

# Full-Mesh Architecture

Virtual Pools for Isolation: Virtual Pool Placement

On each WarpGate:

Say we are to scale up vpool $_i$

For $br_j$ in Bladerunner:

    weight_of_$br_j$ = HASH( vpool$_i$ , $br_j$ )

# Full-Mesh Architecture

Virtual Pools for Isolation: Virtual Pool Placement

On each WarpGate:
Say we are to scale up vpool $_i$
For $br_j$ in Bladerunner:
    weight_of_$br_j$ = HASH( vpool$_i$ , $br_j$ )

Ignore weights of bladerunners that
we don't want to grow into.

Finally, upsize into the bladerunner
host with the highest weight.

# Full-Mesh Architecture

Virtual Pools for Isolation: Virtual Pool Placement

On each WarpGate:
Say we are to scale up vpool $_i$
For br$_j$ in Bladerunner:
    weight_of_br$_j$ = HASH( vpool$_i$ , br$_j$ )

Ignore weights of bladerunners that we don't want to grow into.

Finally, upsize into the bladerunner host with the highest weight.

| Bladerunner 1 | Lorem |
| Bladerunner 2 🔥 | Ipsum |
| Bladerunner 3 | Dolor |
| Bladerunner 4 | Ipsum |
| Bladerunner 5 | |
| Bladerunner 6 | Lorem |
| Bladerunner 7 | Amet |
| ... | |
| Bladerunner 'b' | |

# Full-Mesh Architecture

On each WarpGate:
Say we are to scale up vpool $_i$
For br$_j$ in Bladerunner:
    weight_of_br$_j$ = HASH( vpool$_i$ , br$_j$ )

Ignore weights of bladerunners that we don't want to grow into.

Finally, upsize into the bladerunner host with the highest weight.

| | | |
|---|---|---|
| 17 | **Bladerunner 1** | Lorem |
| 97 | **Bladerunner 2** 🔥 | Ipsum |
| 23 | **Bladerunner 3** | Dolor |
| 89 | **Bladerunner 4** | Ipsum |
| 72 | **Bladerunner 5** | |
| 34 | **Bladerunner 6** | Lorem |
| 29 | **Bladerunner 7** | Amet |
| .. | **...** | |
| 88 | **Bladerunner 'b'** | |

# Full-Mesh Architecture

On each WarpGate:

Say we are to scale up vpool $_i$

For br$_j$ in Bladerunner:

   weight_of_br$_j$ = HASH( vpool$_i$ , br$_j$ )

**<u>Ignore weights of bladerunners that we don't want to grow into.</u>**

Finally, upsize into the bladerunner host with the highest weight.

| | | |
|---|---|---|
| ~~17~~ | **Bladerunner 1** | Lorem |
| ~~97~~ | **Bladerunner 2** 🔥 | Ipsum |
| ~~23~~ | **Bladerunner 3** | Dolor |
| ~~89~~ | **Bladerunner 4** | Ipsum |
| 72 | **Bladerunner 5** | |
| ~~34~~ | **Bladerunner 6** | Lorem |
| ~~29~~ | **Bladerunner 7** | Amet |
| .. | **...** | |
| 88 | **Bladerunner 'b'** | |

# Full-Mesh Architecture

Virtual Pools for Isolation: Virtual Pool Placement

On each WarpGate:
Say we are to scale up $vpool_i$
For $br_j$ in Bladerunner:
$weight\_of\_br_j = HASH(\ vpool_i,\ br_j\ )$

Ignore weights of bladerunners that we don't want to grow into.

Finally, upsize into the bladerunner host with the highest weight.

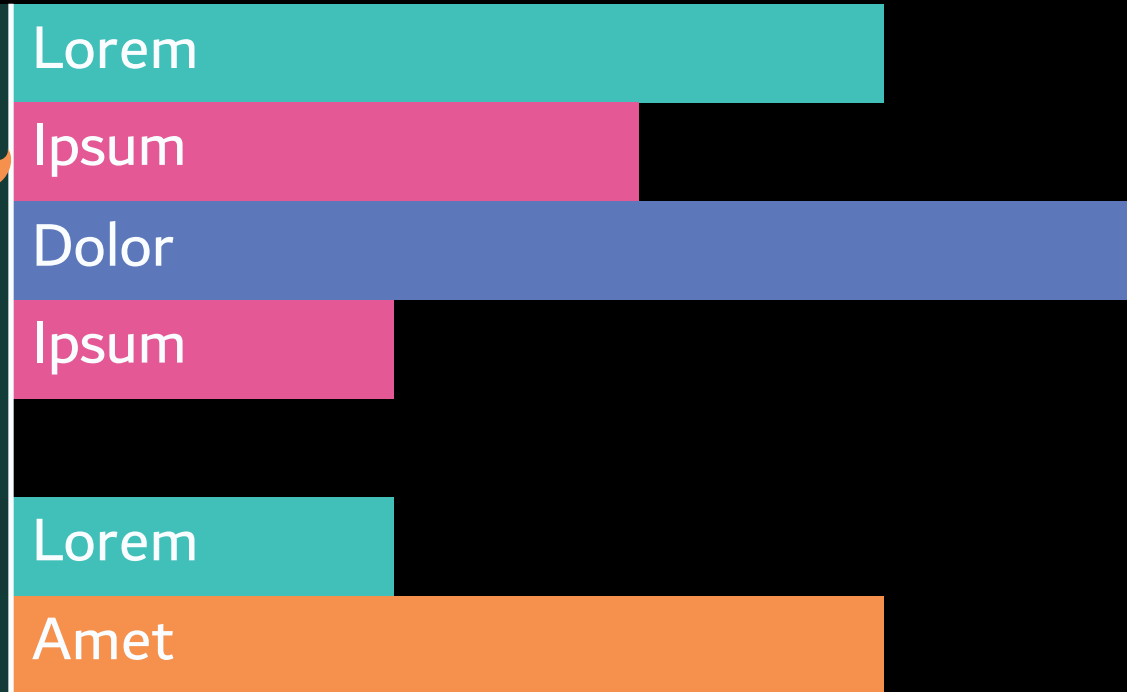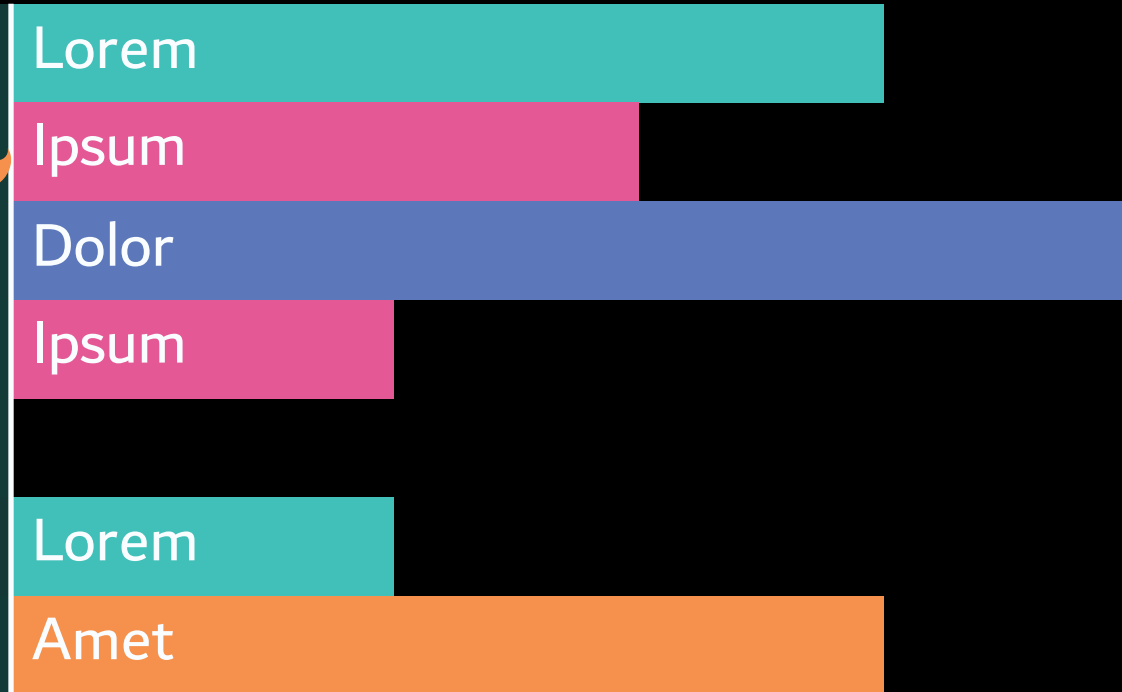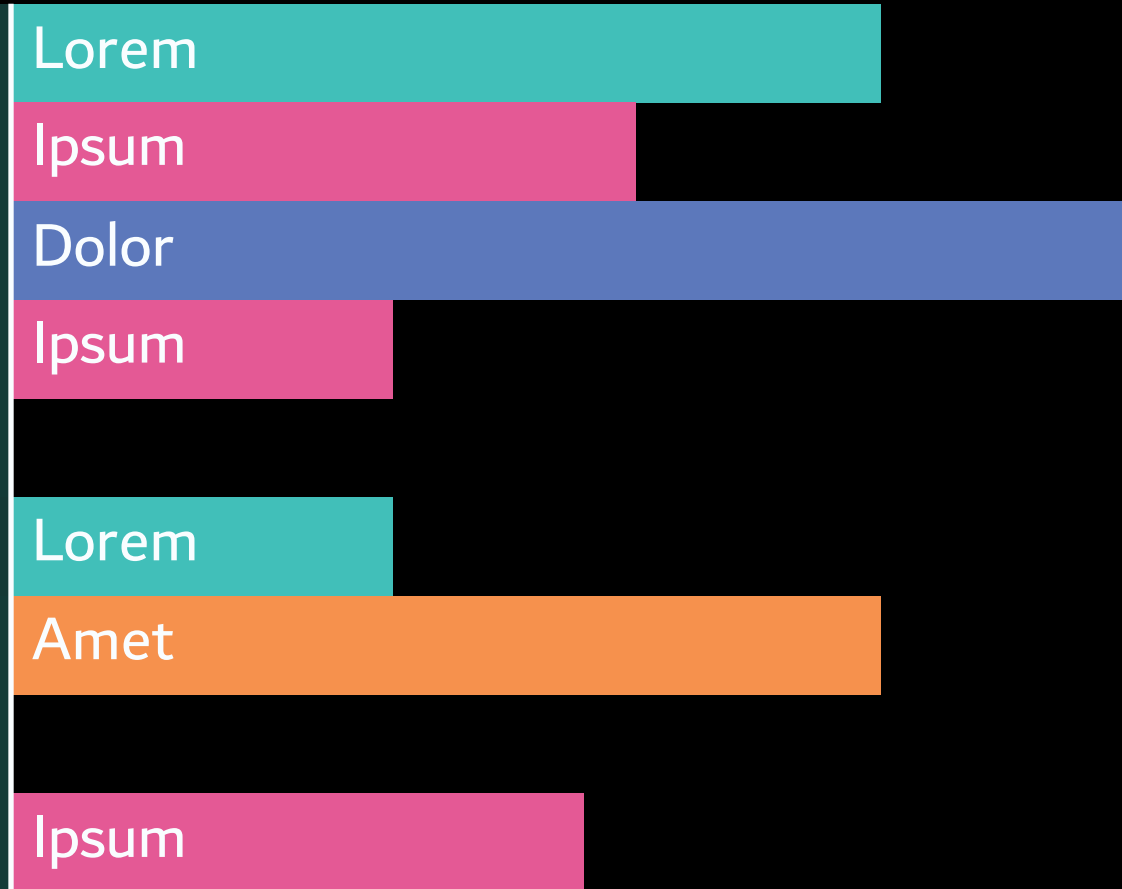| | | |
|---|---|---|
| ~~17~~ | **Bladerunner 1** | Lorem |
| ~~97~~ | **Bladerunner 2** | Ipsum |
| ~~23~~ | **Bladerunner 3** | Dolor |
| ~~89~~ | **Bladerunner 4** | Ipsum |
| 72 | **Bladerunner 5** | |
| ~~34~~ | **Bladerunner 6** | Lorem |
| ~~29~~ | **Bladerunner 7** | Amet |
| .. | **...** | |
| 88 | **Bladerunner 'b'** | Ipsum |

# Full-Mesh Architecture

Virtual Pools for Isolation: Virtual Pool Placement

On each WarpGate:
Say we are to scale down vpool $_i$
For br$_j$ in Bladerunner:
    weight_of_br$_j$ = HASH( vpool$_i$ , br$_j$ )

Ignore weights of bladerunners that don't have desired vpool$_i$

Finally, downsize the bladerunner host with the LOWEST weight in the virtual pool

| Bladerunner 1 | Lorem |
| Bladerunner 2 ❄ | Ipsum |
| Bladerunner 3 | Dolor |
| Bladerunner 4 ❄ | Ipsum |
| Bladerunner 5 | Sit |
| Bladerunner 6 | Lorem |
| Bladerunner 7 | Amet |
| ... | |
| Bladerunner 'b' | Ipsum |

# Full-Mesh Architecture

On each WarpGate:

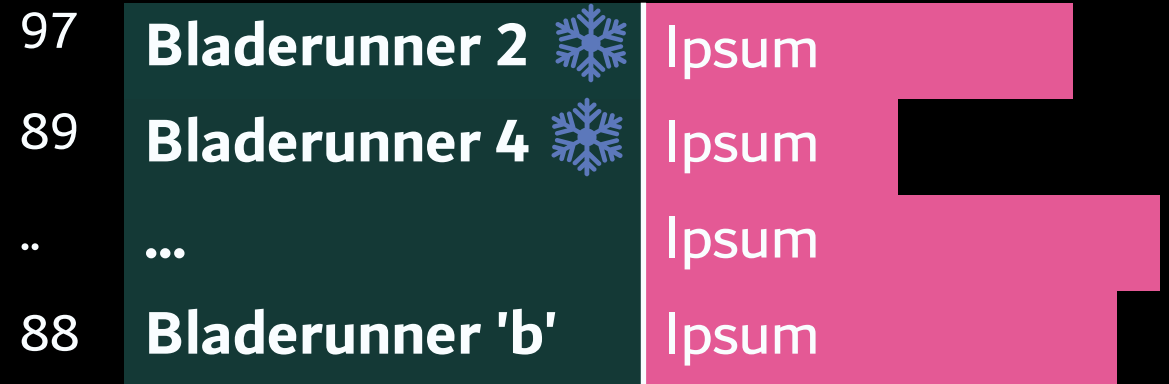Say we are to scale down vpool $_i$

For br$_j$ in Bladerunner:

    weight_of_br$_j$ = HASH( vpool$_i$ , br$_j$ )

**Ignore weights of bladerunners that don't have desired vpool$_i$**

Finally, downsize the bladerunner host with the LOWEST weight in the virtual pool

## Virtual Pool of Ipsum hosts

| | | |
|---|---|---|
| 97 | **Bladerunner 2** ❄️ | Ipsum |
| 89 | **Bladerunner 4** ❄️ | Ipsum |
| .. | **...** | Ipsum |
| 88 | **Bladerunner 'b'** | Ipsum |

# Full-Mesh Architecture

On each WarpGate:

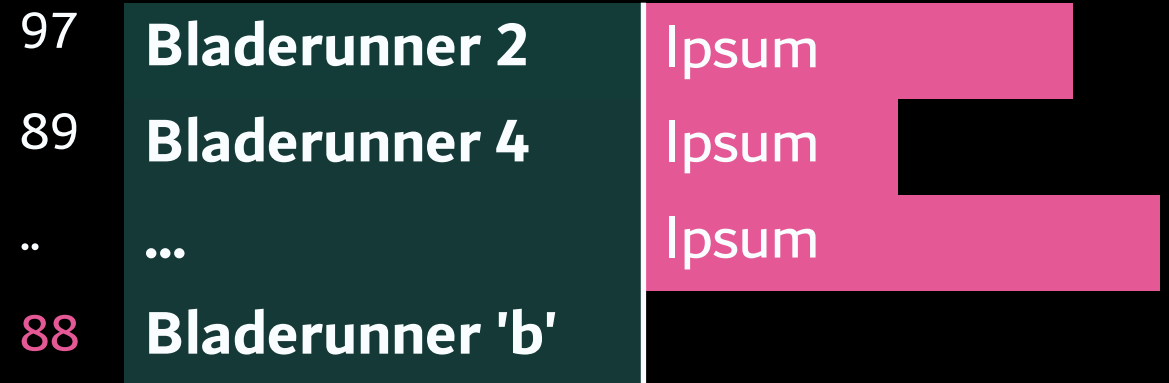Say we are to scale down vpool $_i$

For br$_j$ in Bladerunner:

   weight_of_br$_j$ = HASH( vpool$_i$ , br$_j$ )

**<u>Ignore weights of bladerunners that</u>**
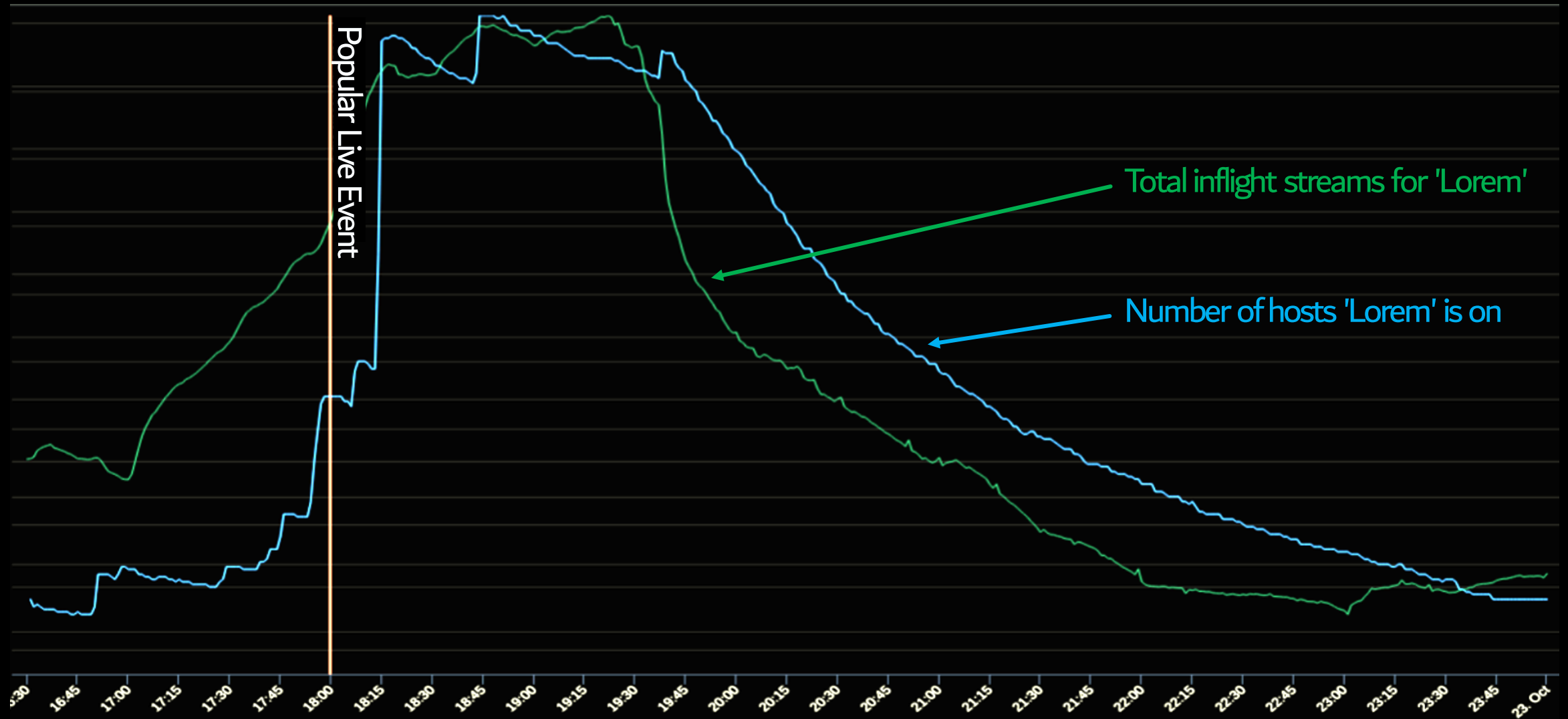**<u>don't have desired  vpool$_i$</u>**

Finally, downsize the bladerunner
host with the LOWEST weight in the
virtual pool

## Virtual Pool of Ipsum hosts

| | | |
|---|---|---|
| 97 | **Bladerunner 2** | Ipsum |
| 89 | **Bladerunner 4** | Ipsum |
| .. | **...** | Ipsum |
| 88 | **Bladerunner 'b'** | |

# Full-Mesh Architecture

Virtual Pools for Isolation: Case Study of a Recent Event

# Full-Mesh Architecture

Virtual Pools for Isolation

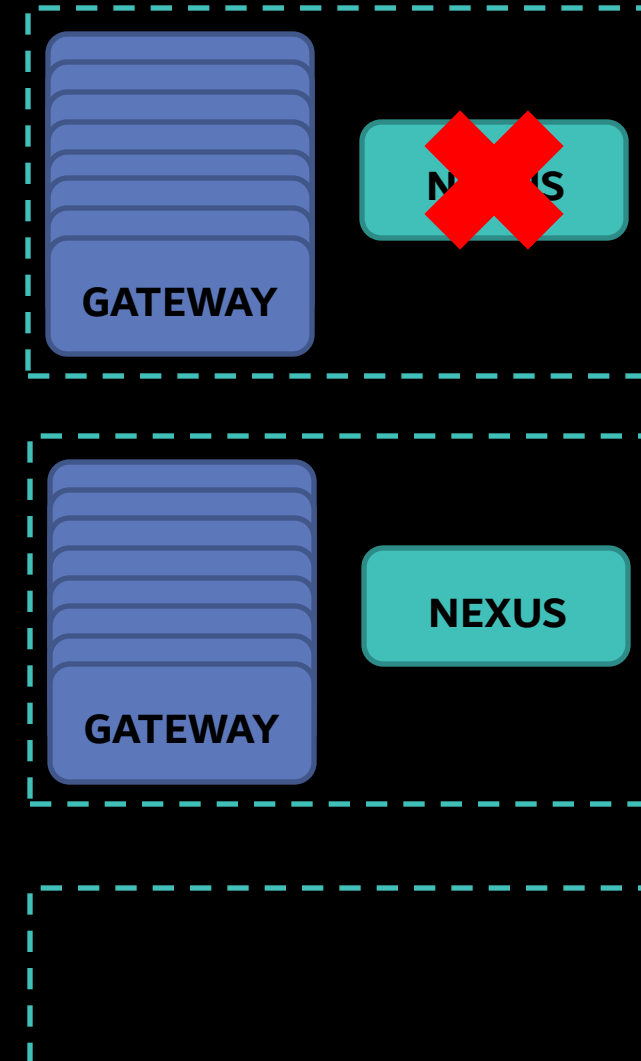Noisy Neighbour

Resource Attribution

Maintenance Costs

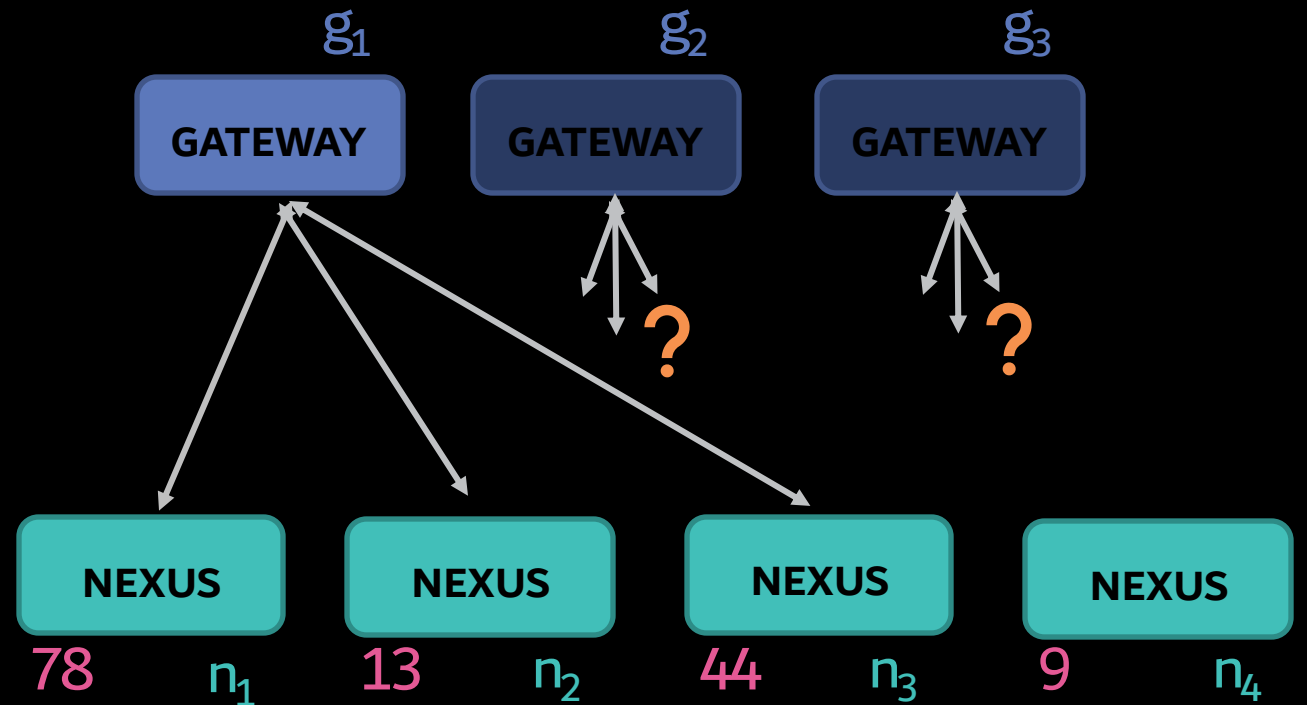# The Evolution of Traffic Routing in a Streaming World

1. Moving away from Collocation

# The Evolution of Traffic Routing in a Streaming World

1. Moving away from Collocation

2. Rendezvous Hashing

# The Evolution of Traffic Routing in a Streaming World

1. Moving away from Collocation

2. Rendezvous Hashing

3. Full Mesh Architecture

| | |
|---|---|
| **Bladerunner 1** | Lorem |
| **Bladerunner 2** | Ipsum |
| **Bladerunner 3** | Dolor |
| **Bladerunner 4** | Ipsum |
| **Bladerunner 5** | Sit |
| **Bladerunner 6** | Lorem |
| **Bladerunner 7** | Amet |
| **...** | |
| **Bladerunner 'b'** | Dolor |

# Thank you