

KIRA – Scalable Zero-Touch Routing

Roland Bless, Martina Zitterbart
Institute of Telematics, KIT

Zoran Despotovic, Artur Hecker
Huawei Research Center, Munich

KIRA: Kademlia-directed ID-based Routing Architecture

R. Bless, M. Zitterbart, Z. Despotovic and A. Hecker, „KIRA: Distributed Scalable ID-based Routing with Fast Forwarding“, 2022 IFIP Networking Conference (IFIP Networking), 2022, <https://s.kit.edu/KIRA>

Public Side Meeting – IETF Rules Apply

- All public side meetings are subject to the [IETF Note Well](#)
- Participants are expected to follow the usual IETF policies on personal conduct, IPR disclosure obligations, etc.
- You are contributing to IETF work
- Policies on patents and code of conduct apply
- Disclose relevant IPR or do not make contributions
- Be respectful and courteous even/especially when you disagree

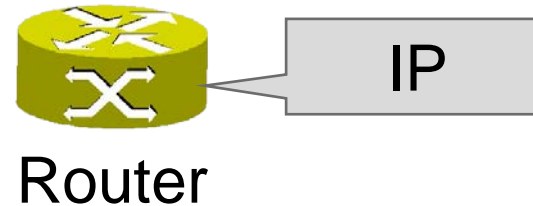
Agenda

- Introduction to KIRA (20min)
- Questions & Answers
- Discuss next steps
 - Supporting the IETF activities

Introduction

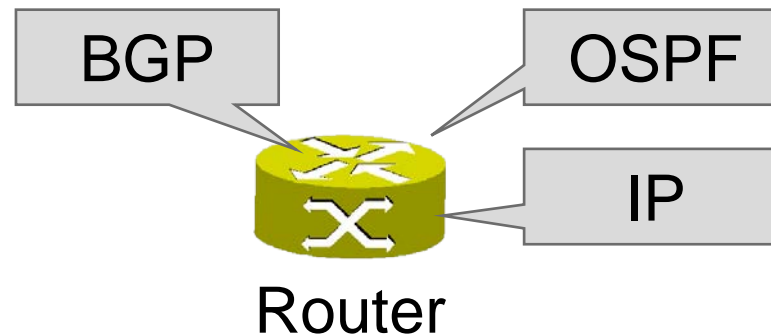
Internet Infrastructure...

- is becoming more complex
 - higher interdependencies of services
- must be reliable → resilient operation



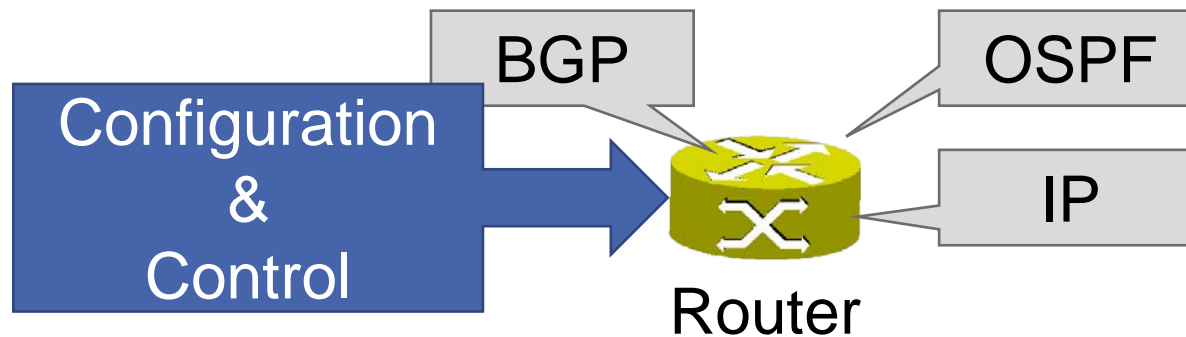
Internet Infrastructure...

- is becoming more complex
 - higher interdependencies of services
- must be reliable → resilient operation



Internet Infrastructure...

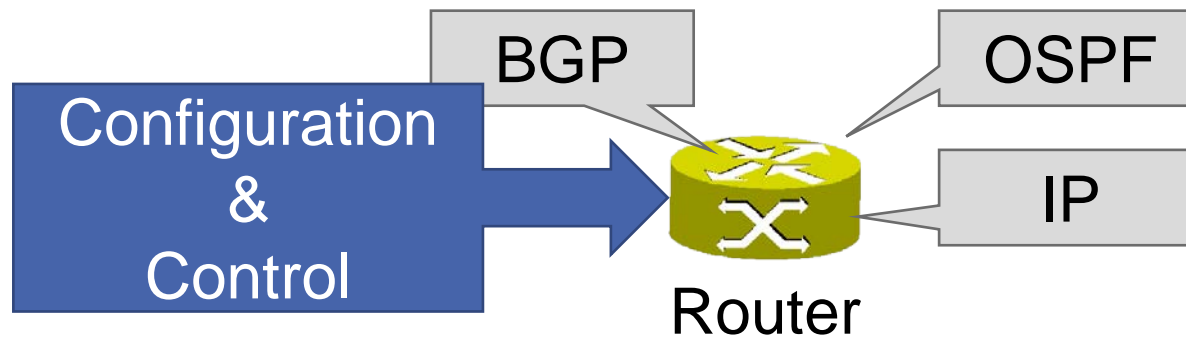
- is becoming more complex
 - higher interdependencies of services
- must be reliable → resilient operation



- Requires configuration via management/control plane

Internet Infrastructure...

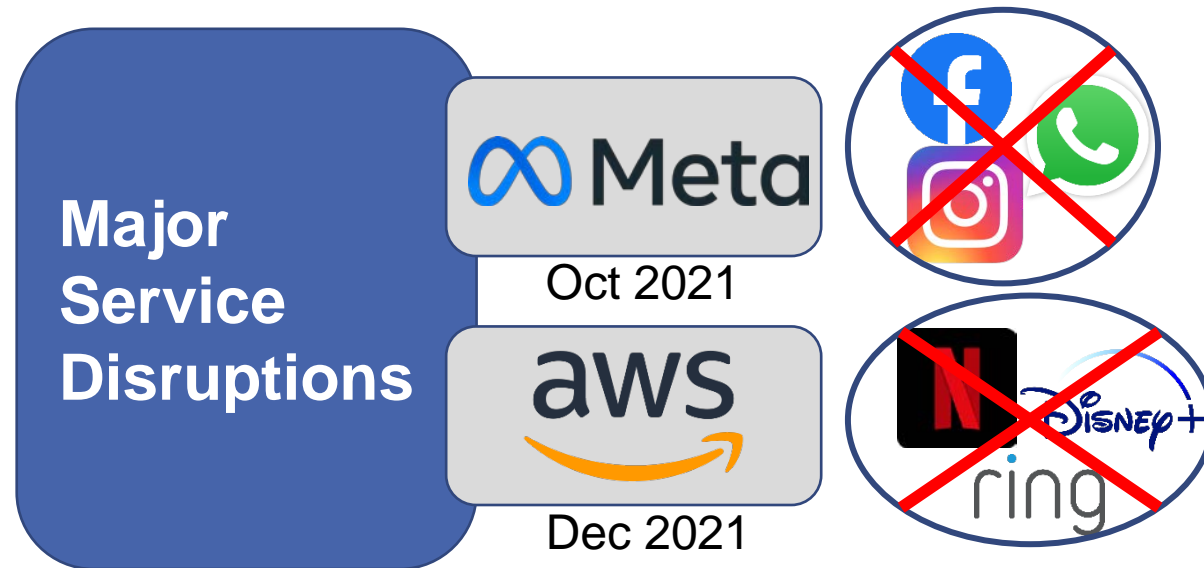
- is becoming **more complex**
 - higher **interdependencies** of services
- must be reliable → **resilient operation**



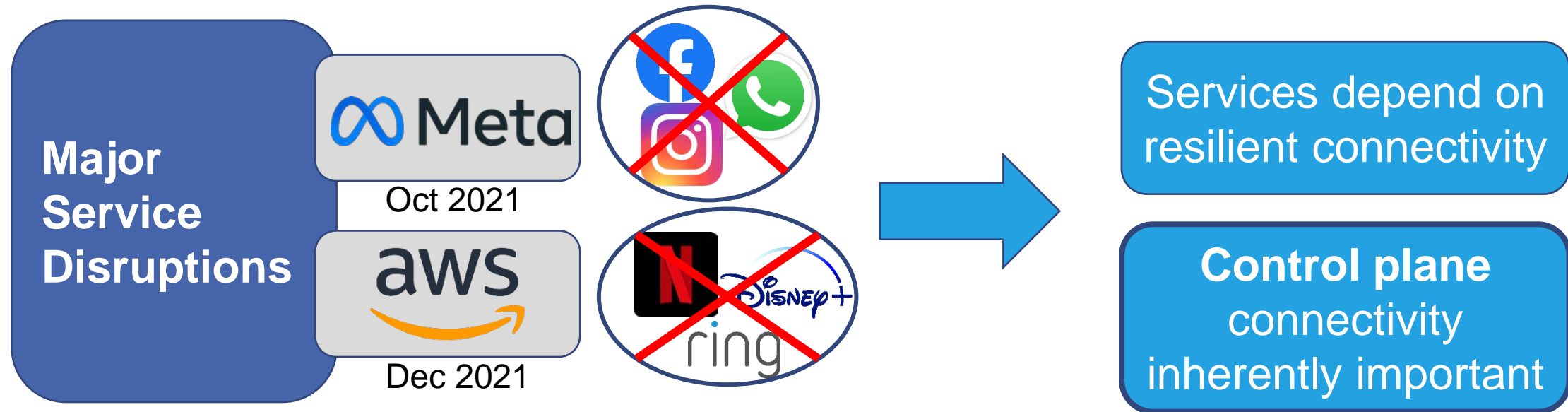
- Requires configuration via **management/control plane**

Working Control Plane Connectivity
=
Foundation for Resilient Internet Infrastructures

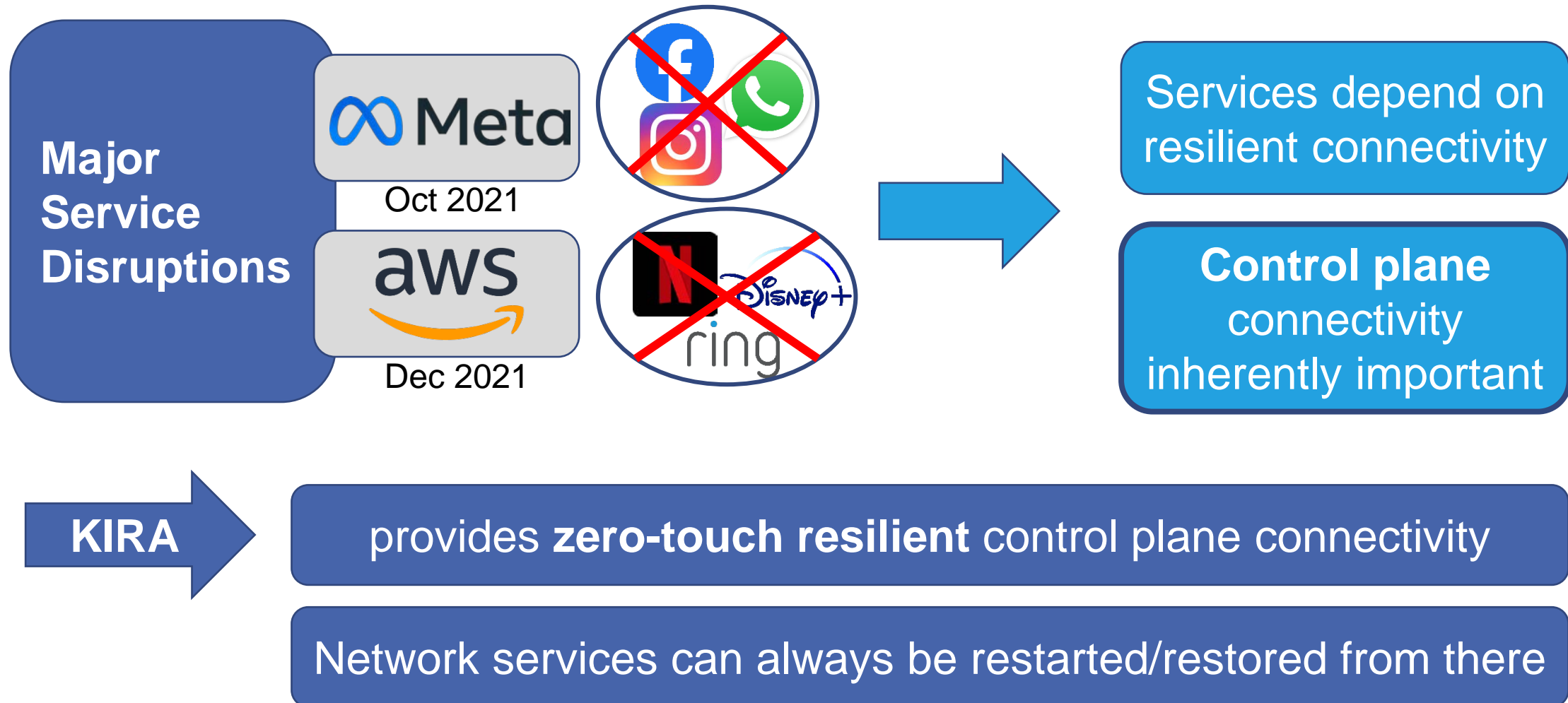
Controllability and Control Planes



Controllability and Control Planes



Controllability and Control Planes



Infrastructure Complexity

A Decentralized SDN Architecture for the WAN

Alexander Krentzel
Google / UC Berkeley

Nitika Saran*
Cornell

Bikash Koley
Google

Subhasree Mandal
Google

Ashok Narayanan
Google

Sylvia Ratnasamy
Google / UC Berkeley

Ali Al-Shabibi
Google

Anees Shaikh
Google

Rob Shakir
Google

Ankit Singla
Google

Hakim Weatherspoon
Cornell

dsdn-sigcomm@google.com

SIGCOMM 2024

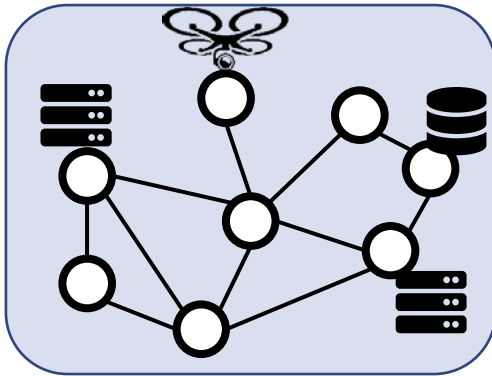
The fundamental challenge in eliminating major failures lies in their **complexity**: they typically involve bugs or errors in multiple components spanning diverse teams and codebases, that **interact in unanticipated ways**. As a result, these outages persist, despite our extensive efforts to improve testing, diagnostics, change procedures, and verification.



...address complexity directly and focus on simplifying the network.

Control Planes of Future Networks Need to Support...

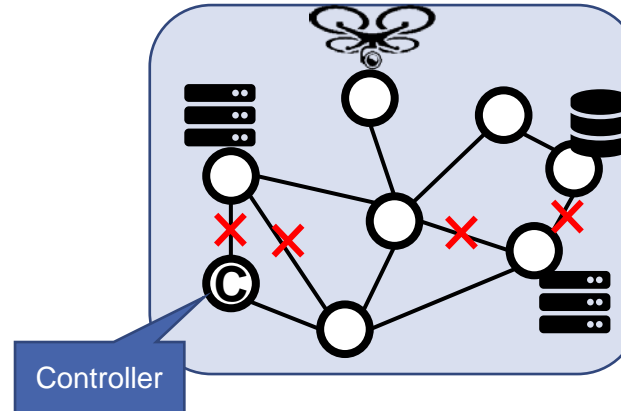
Interconnection of a Large Pool of Networked Resources



Compute, Storage, Network

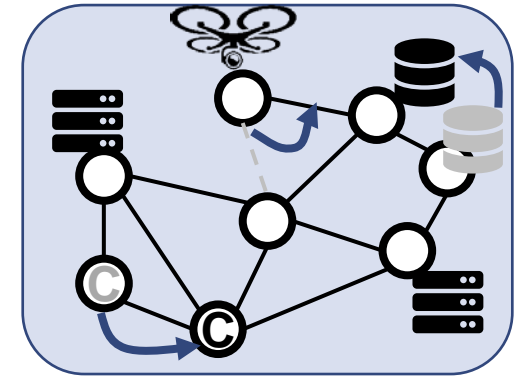
- Scalability
- In-band control
- High dynamics
- Multiple domains
- Various topologies

Resilient Connectivity for Control Plane



- Zero-touch
- Fast convergence
- Network split
- Nomadic networks

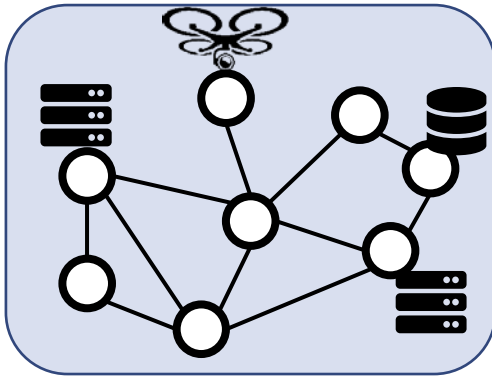
Stable Addresses for Moving Resources



- ID-based addresses

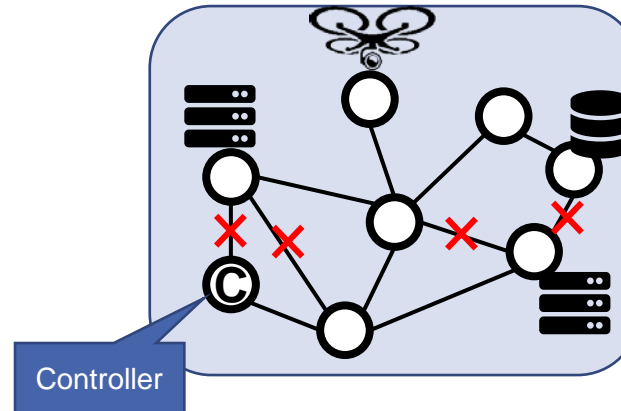
What KIRA achieves...

Interconnection of a Large Pool of Networked Resources

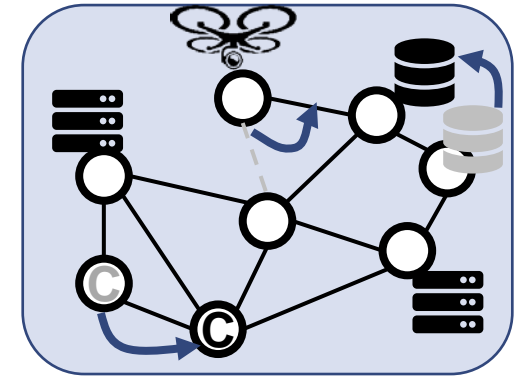


Compute, Storage, Network

Resilient Connectivity for Control Plane



Stable Addresses for Moving Resources



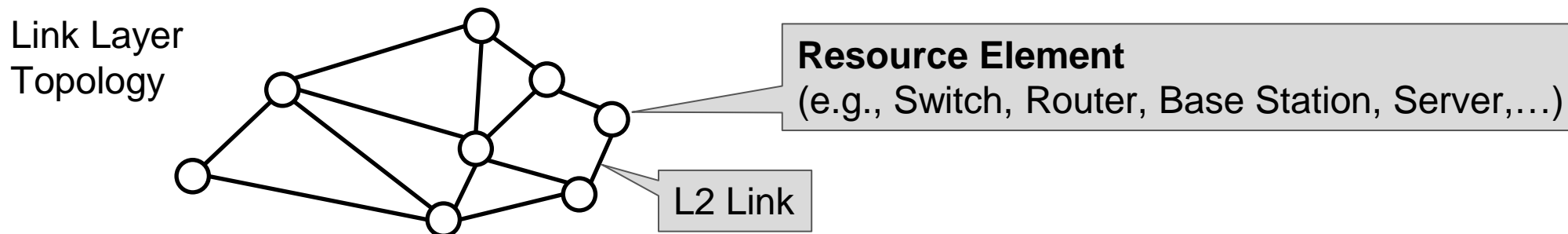
■ KIRA provides (all-in-one)

- Massive scalability (100,000s of nodes)
- Zero-touch (no configuration + adaptation)
- Dynamics: fast convergence, loop free, fast reroute
- Topological versatility
- Efficient routes

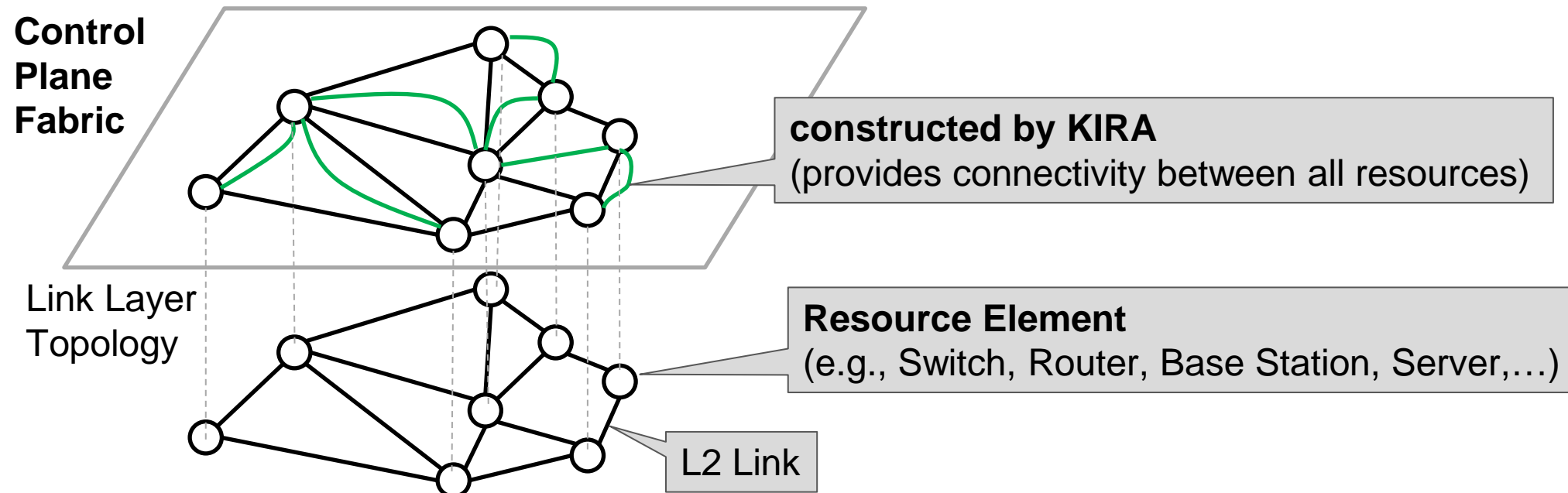
■ Related Works (examples)

- UIP: lacks dynamics and efficient routes
- DISCO: lacks dynamics
- RIFT, Data Center BGP/OSPF/IS-IS: specific topologies only, not ID-based
- RPL: traffic concentration near root, zero-touch?
- ...

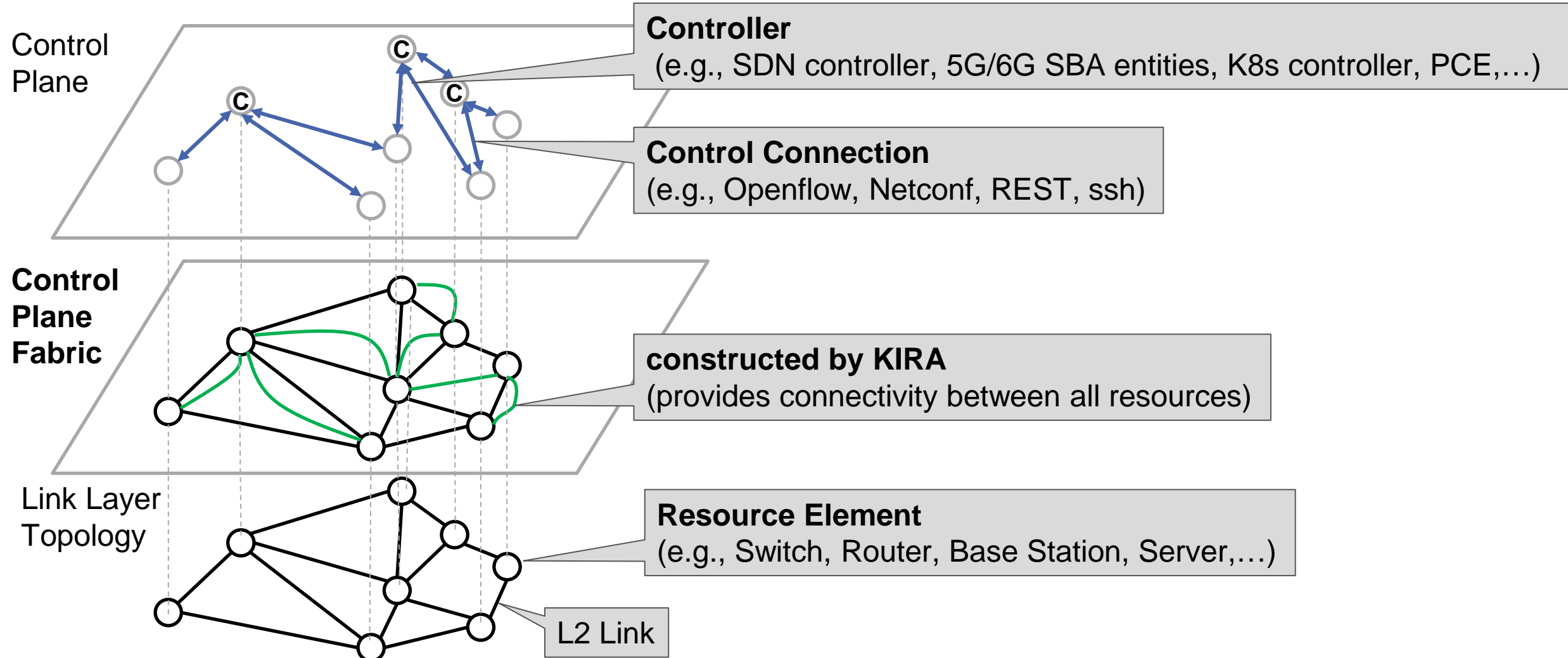
What KIRA provides...



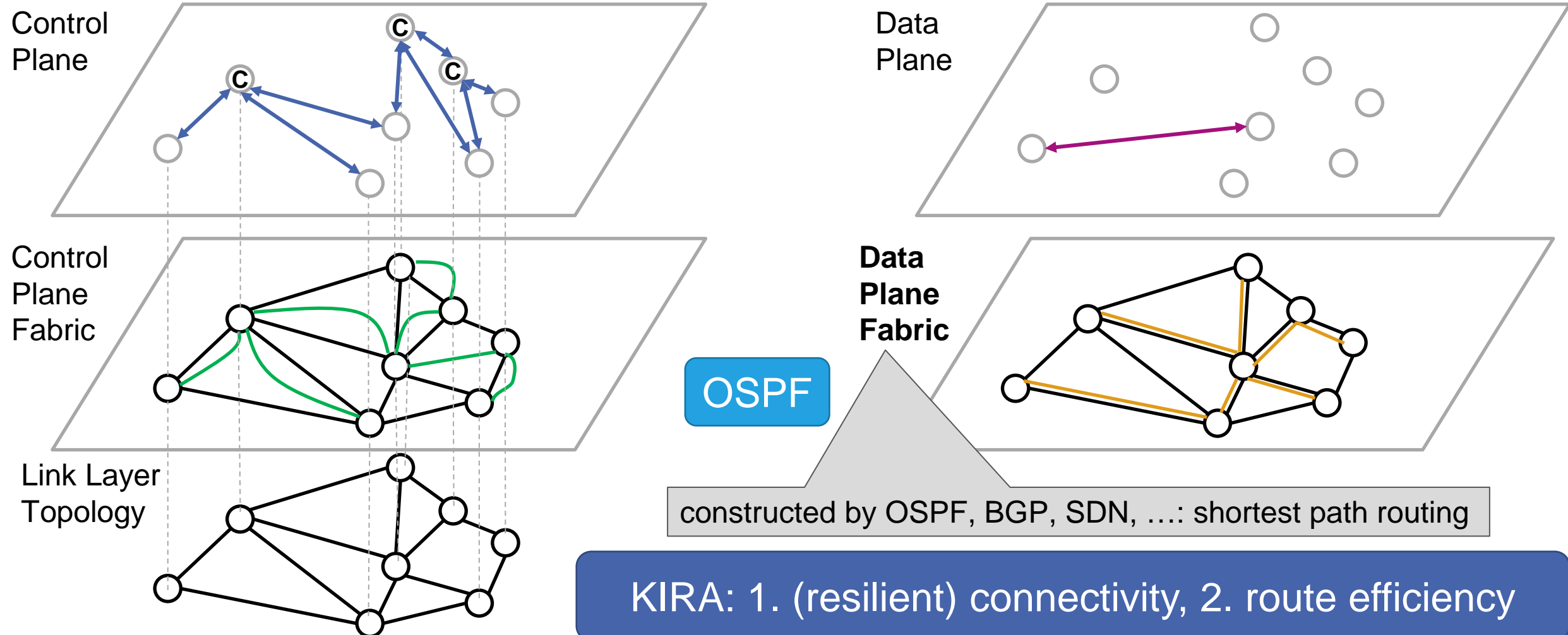
What KIRA provides...



What KIRA provides...

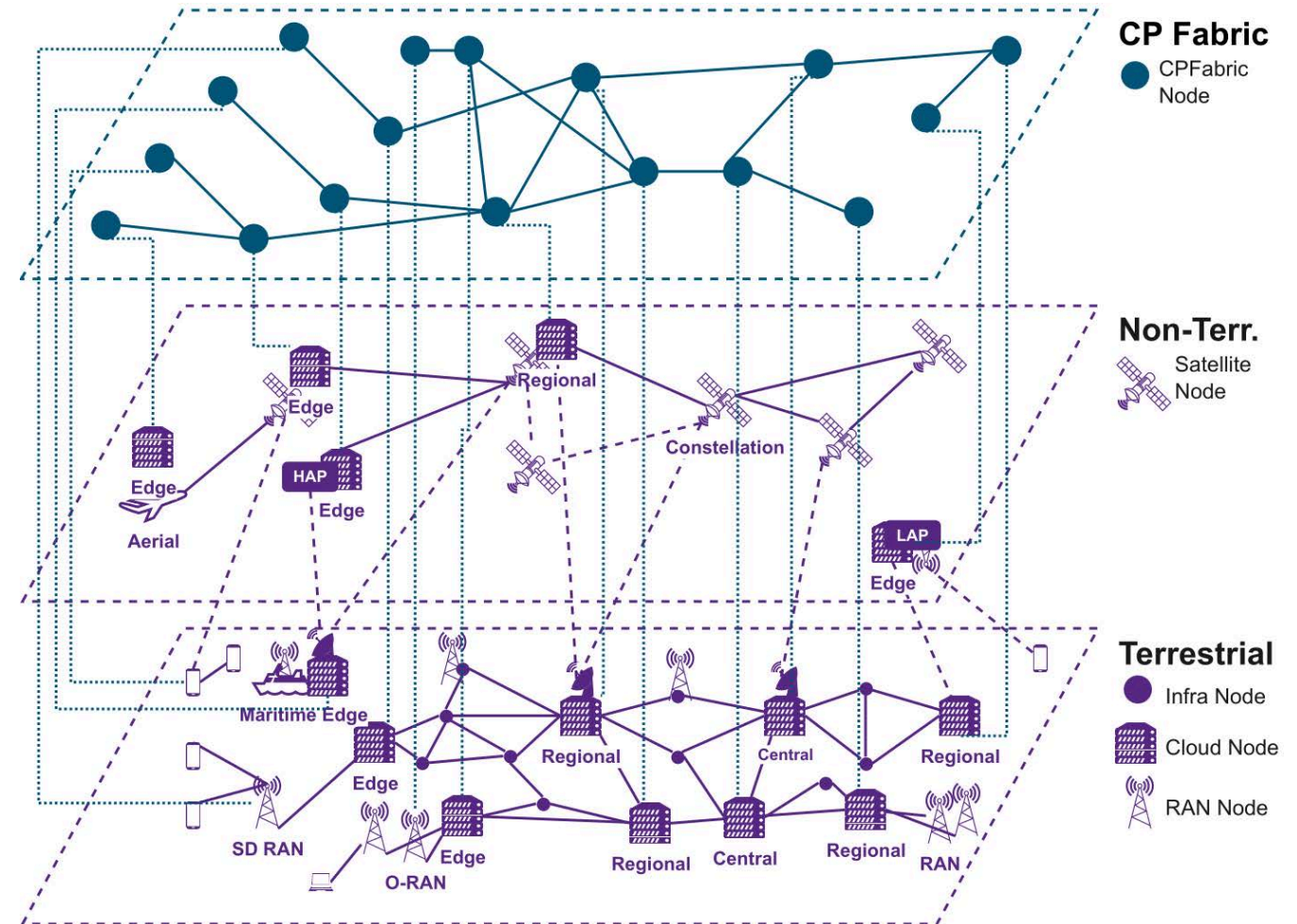


What KIRA provides...



Use Case – 6G Control Plane

- Non-terrestrial Networks (Drones, Satellites)
→ dynamic and **mobile**
- **Nomadic Networks**
→ autonomous, self-organizing control plane
- 10^6 of base stations in China in a single provider network
→ scalability
- Single routing solution that interconnects everything for the Control Plane



Example – Networking Control and Management

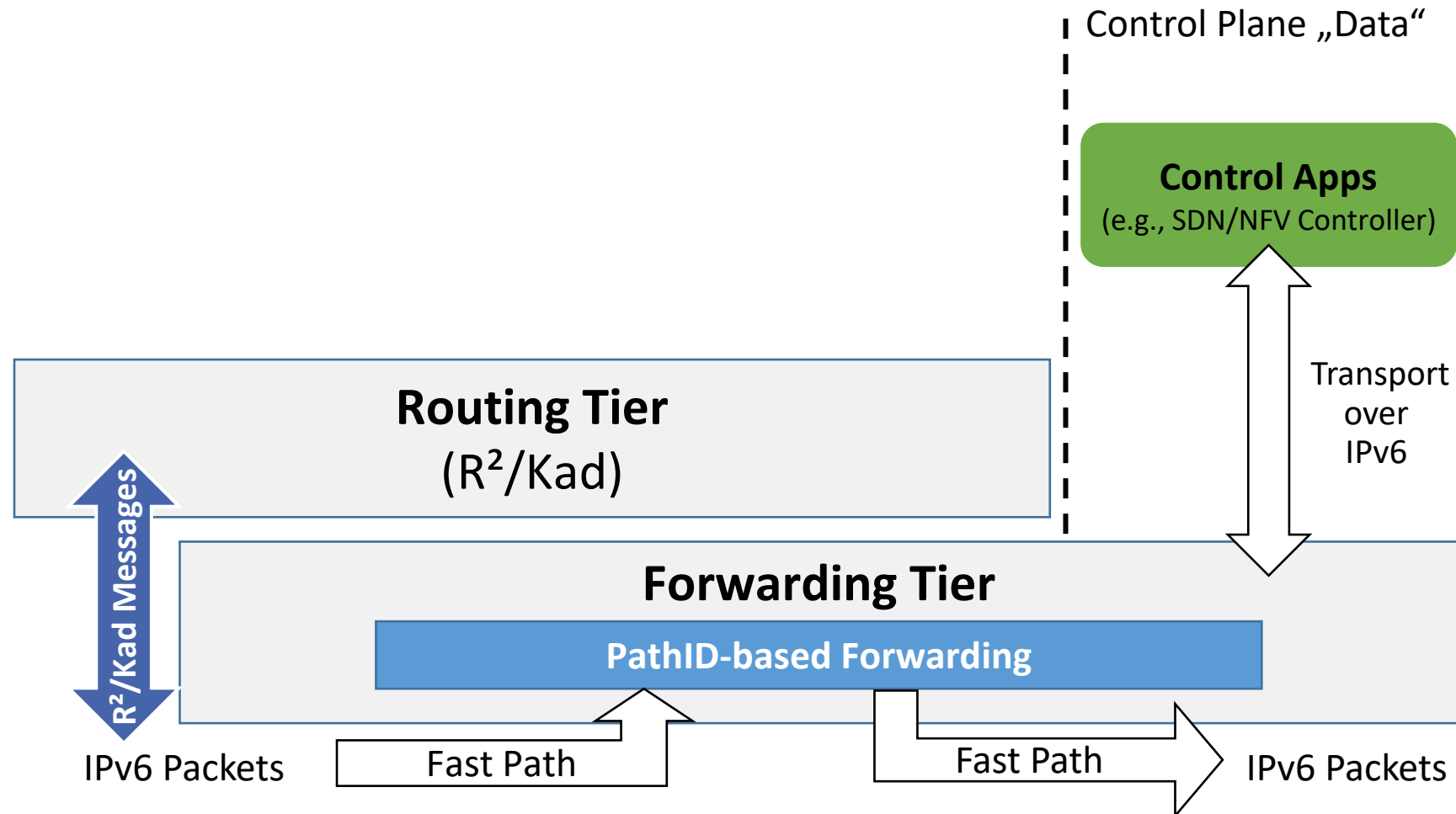
- KIRA bootstraps the control plane fabric
- Resources can register themselves to be found (e.g., run ANIMA on top)
- Topologically dependent routing can be built on top
 - e.g., divide the network into areas, assign and distribute prefixes, configure OSPF routers etc.
- Distributed controllers may claim control over some resource subset
- In case things go (horribly) wrong
 - restart from scratch, revert to last working configuration

Never lose control over your network!

Features

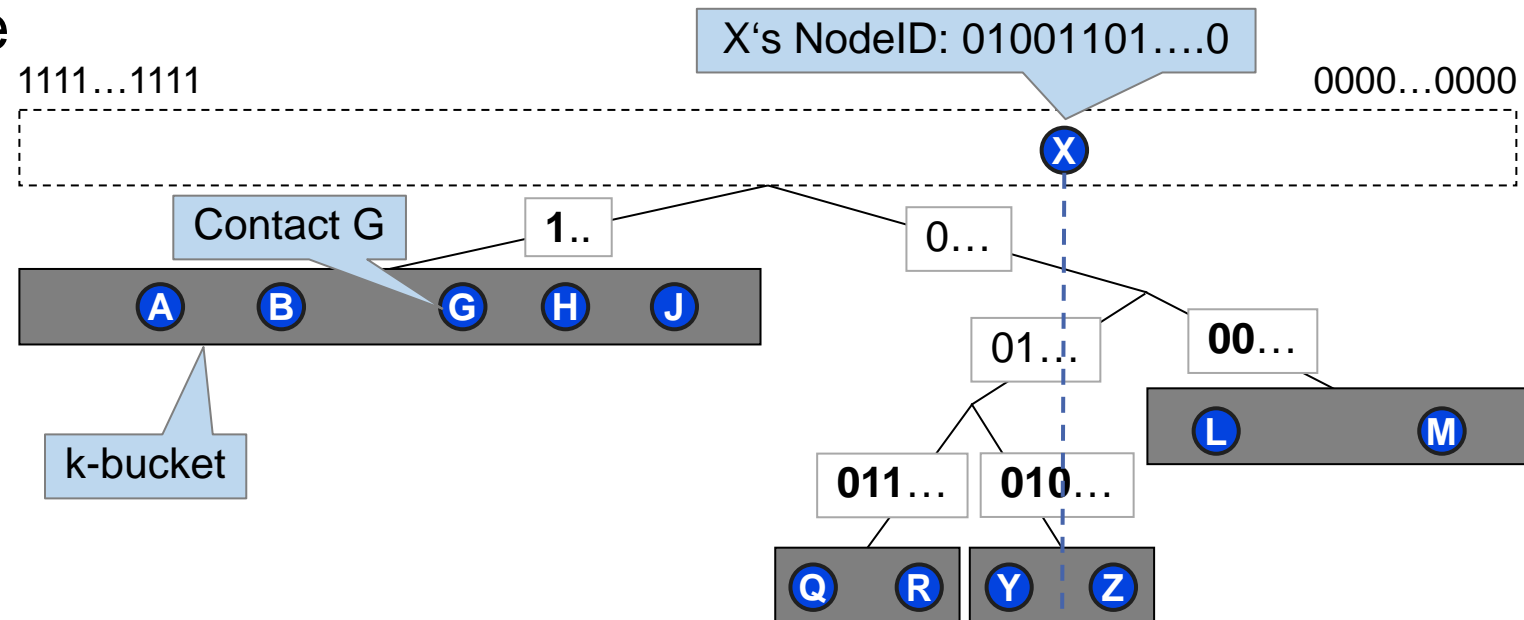
- **Kademlia as ID-based Overlay (XOR metric)**
 - 112bit NodeID as address → IPv6 address w/ 16bit prefix
 - Similar to Virtual Ring Routing (that is less efficient)
 - **Small routing tables** – $O(\log n)$ n : number of nodes → Stretch as trade-off
 - Supports multi-homing and mobility
- **Proximity Neighbor Selection and Proximity Routing**
- **Path Discovery** (First Path, Later Path) + **Path Rediscovery** (Dynamics)
- **Loop-free** even during convergence
- Source routing for routing messages + path validation
- **PathID-based forwarding for data**
 - no special forwarding hardware required (IPv6 support)

Architecture



Routing Tier – Routing Table

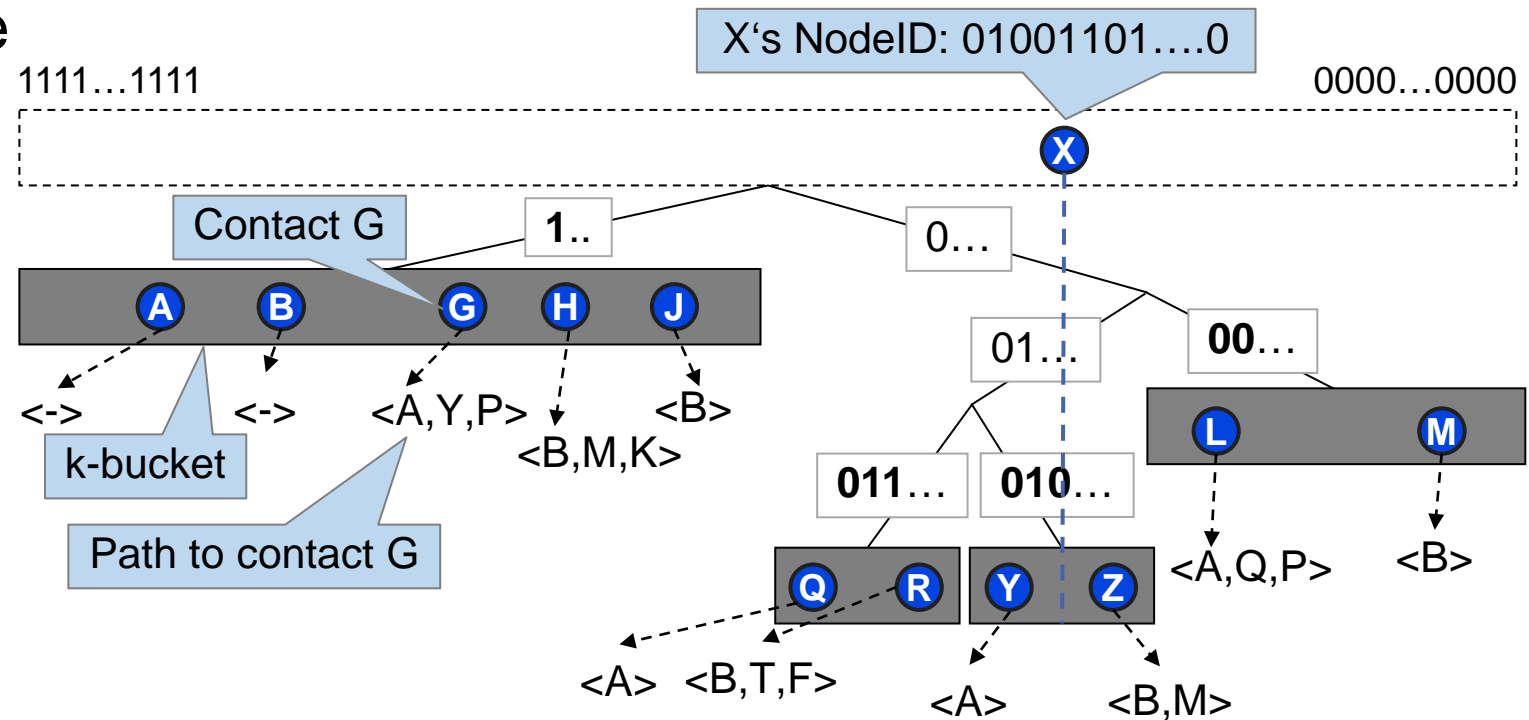
- Tree of **buckets** holding up to **k contacts** (e.g., $k=40$)
 - Arranged by XOR distance
 - Bucket split if it contains own NodeID



Routing table size: $O(l_G \log n)$, l_G average path length

Routing Tier – Routing Table

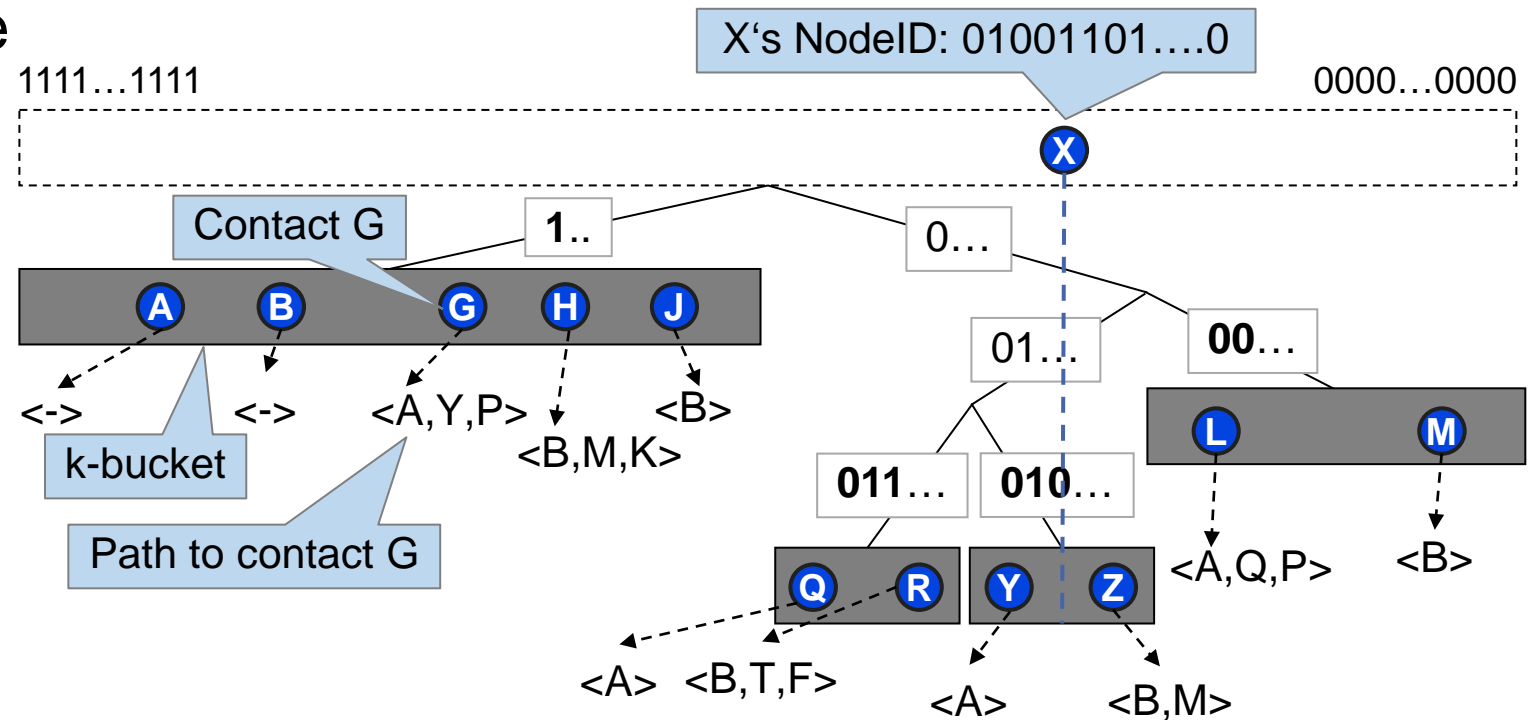
- Tree of **buckets** holding up to **k contacts** (e.g., $k=40$)
 - Arranged by XOR distance
 - Bucket split if it contains own NodeID
- **Path vectors** are stored for each contact



Routing table size: $O(l_G \log n)$, l_G average path length

Routing Tier – Routing Table

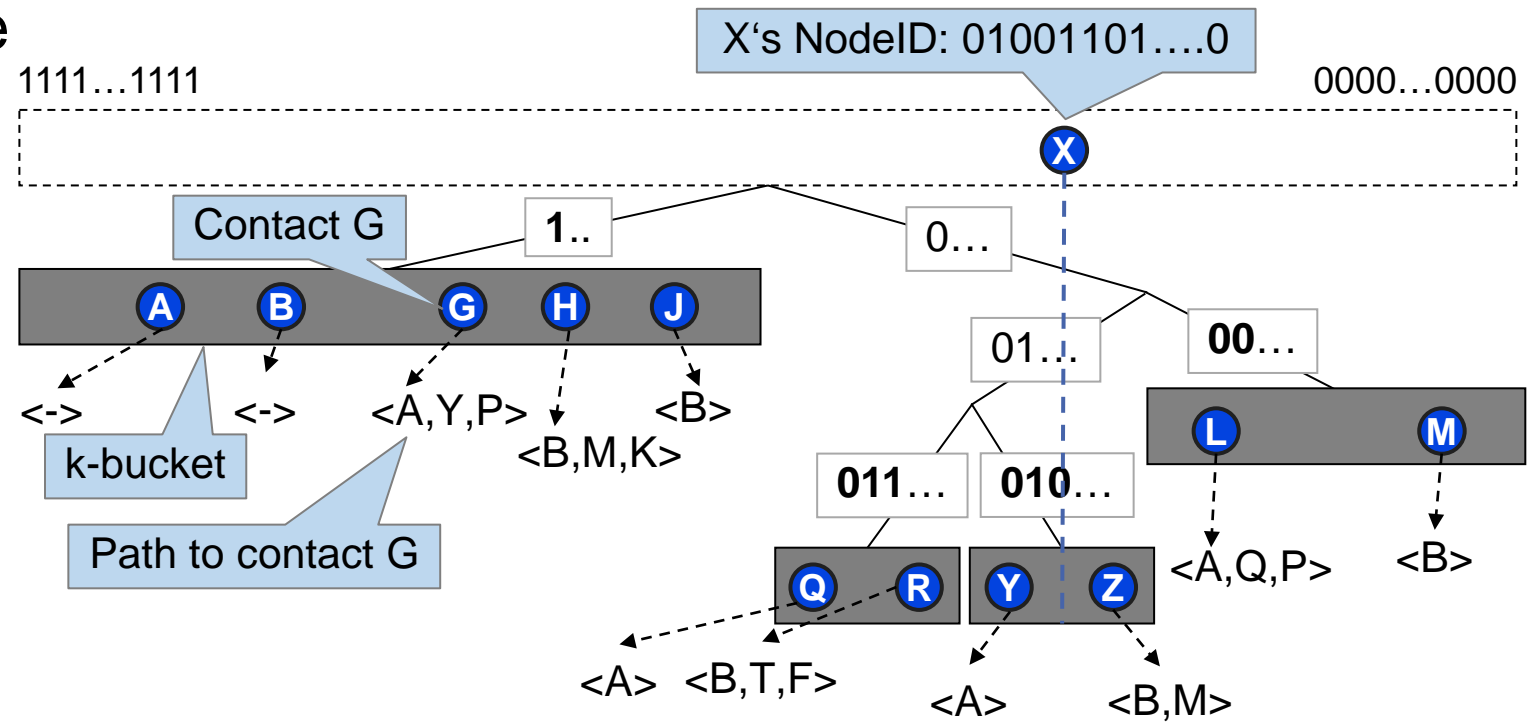
- Tree of **buckets** holding up to **k contacts** (e.g., $k=40$)
 - Arranged by XOR distance
 - Bucket split if it contains own NodeID
- **Path vectors** are stored for each contact
- **Efficient routes**
 - Shorter routes preferred



Routing table size: $O(l_G \log n)$, l_G average path length

Routing Tier – Routing Table

- Tree of **buckets** holding up to **k contacts** (e.g., $k=40$)
 - Arranged by XOR distance
 - Bucket split if it contains own NodeID
- **Path vectors** are stored for each contact
- **Efficient routes**
 - Shorter routes preferred
- Size k of k -buckets can be set per node
→ **flexible memory / stretch trade-off**



Routing table size: $O(l_G \log n)$, l_G average path length

Routing Tier – Routing Table

- Tree of **buckets** holding up to **k contacts** (e.g., $k=40$)

- Arranged by XOR distance
- Bucket split if it contains own NodeID

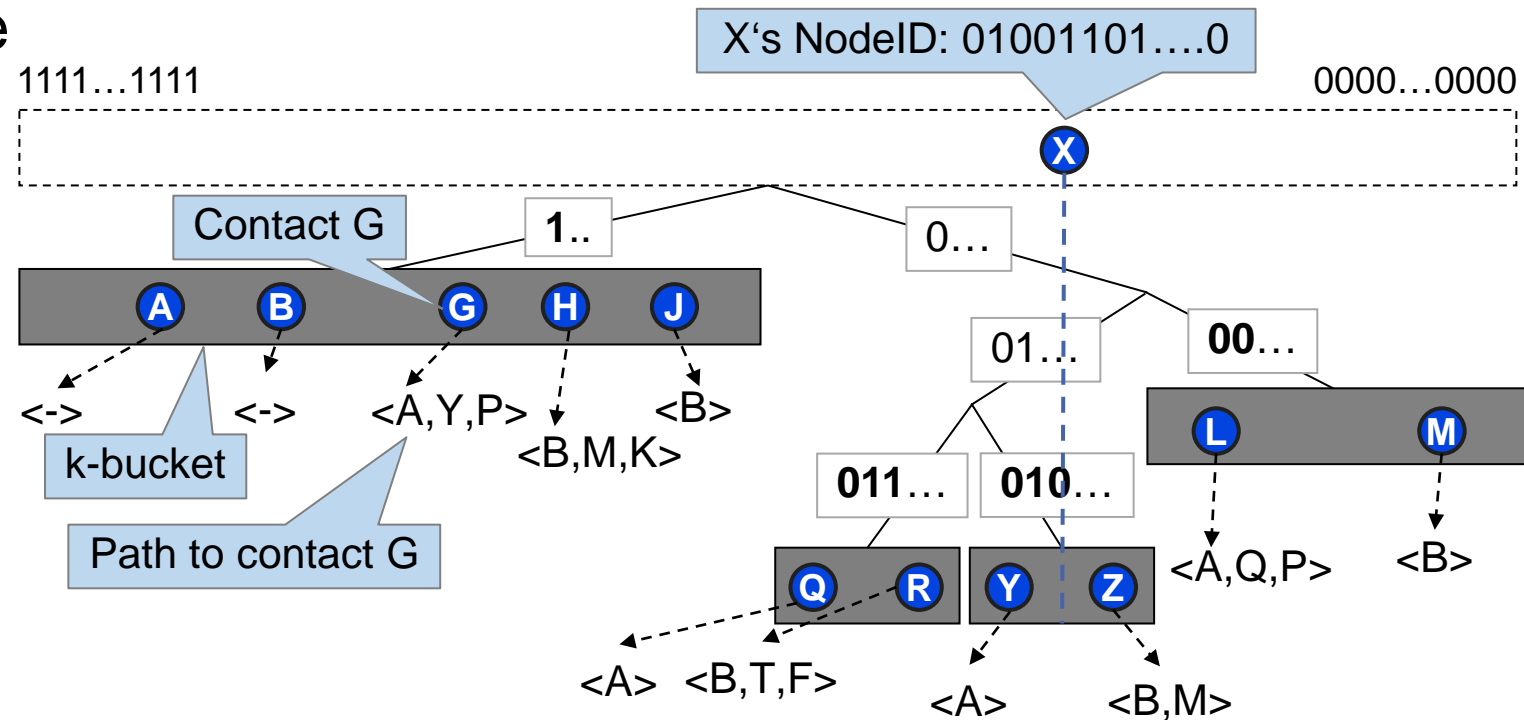
- Path vectors** are stored for each contact

- Efficient routes**

- Shorter routes preferred

- Size k of k -buckets can be set per node
→ **flexible memory / stretch trade-off**

- Each node routes to XOR-wise closest entry



Routing table size: $O(l_G \log n)$, l_G average path length

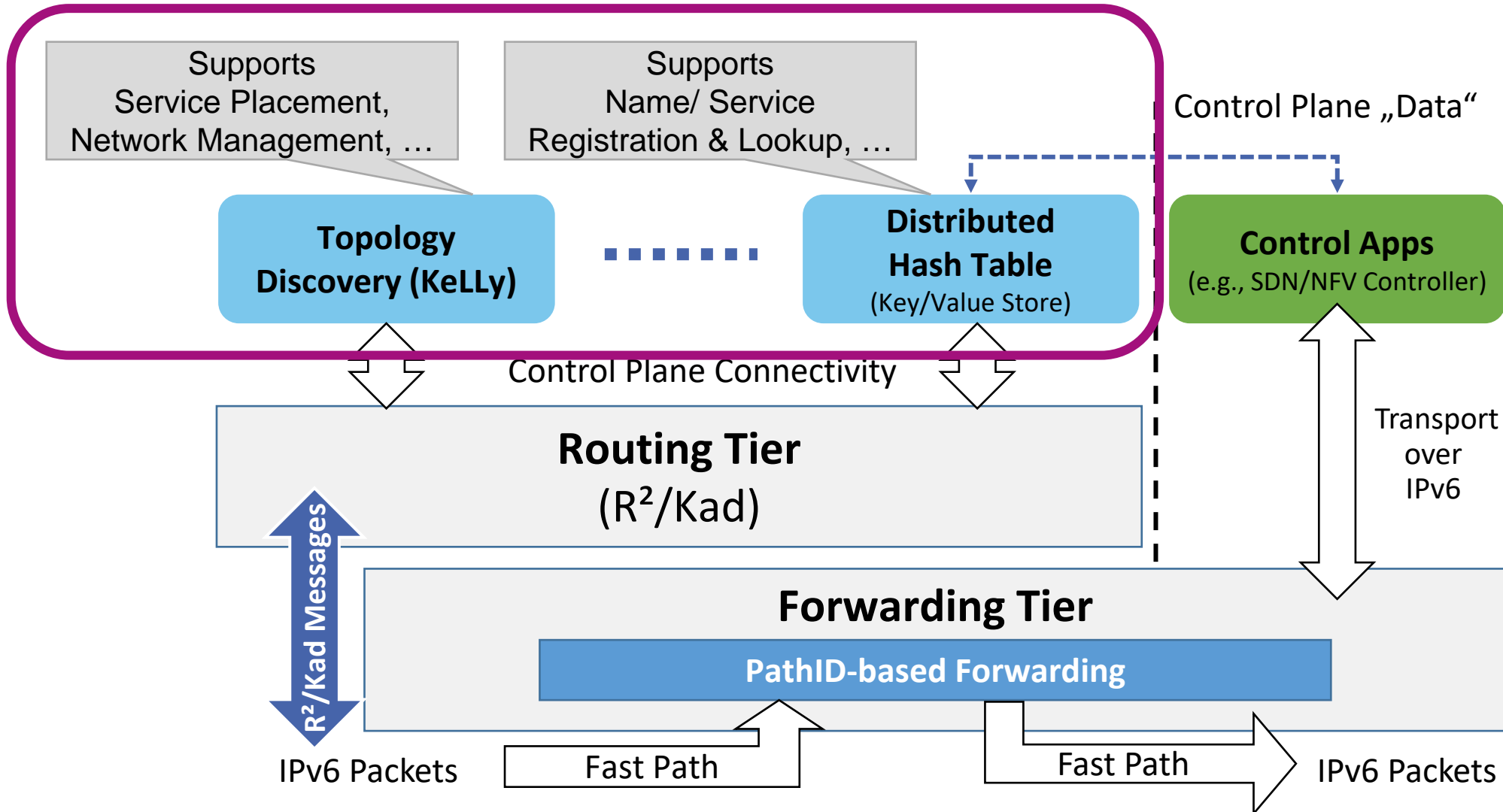
Forwarding Tier

- Approach: replace source routes with PathIDs
 - $\text{PathID}(\langle A, Q, M, Z \rangle) = \text{Hash}(A \mid Q \mid M \mid Z)$
- Need PathID in addition to NodeIDs → encapsulation (e.g., SRv6, GRE)
- Use PathID as unique label for path segment → Label Switching
- Precalculate PathIDs for 2-hop (physical) vicinity
- Explicit path setup for paths ≥ 6 hops
- Cost?
 - 100,000 nodes powerlaw topology:
99th percentile forwarding table ~ 2500 entries
- Uses existing forwarding plane mechanisms
- Works with Longest Matching Prefix, nftables, SDN, eBPF etc.

Further Features (Continued)

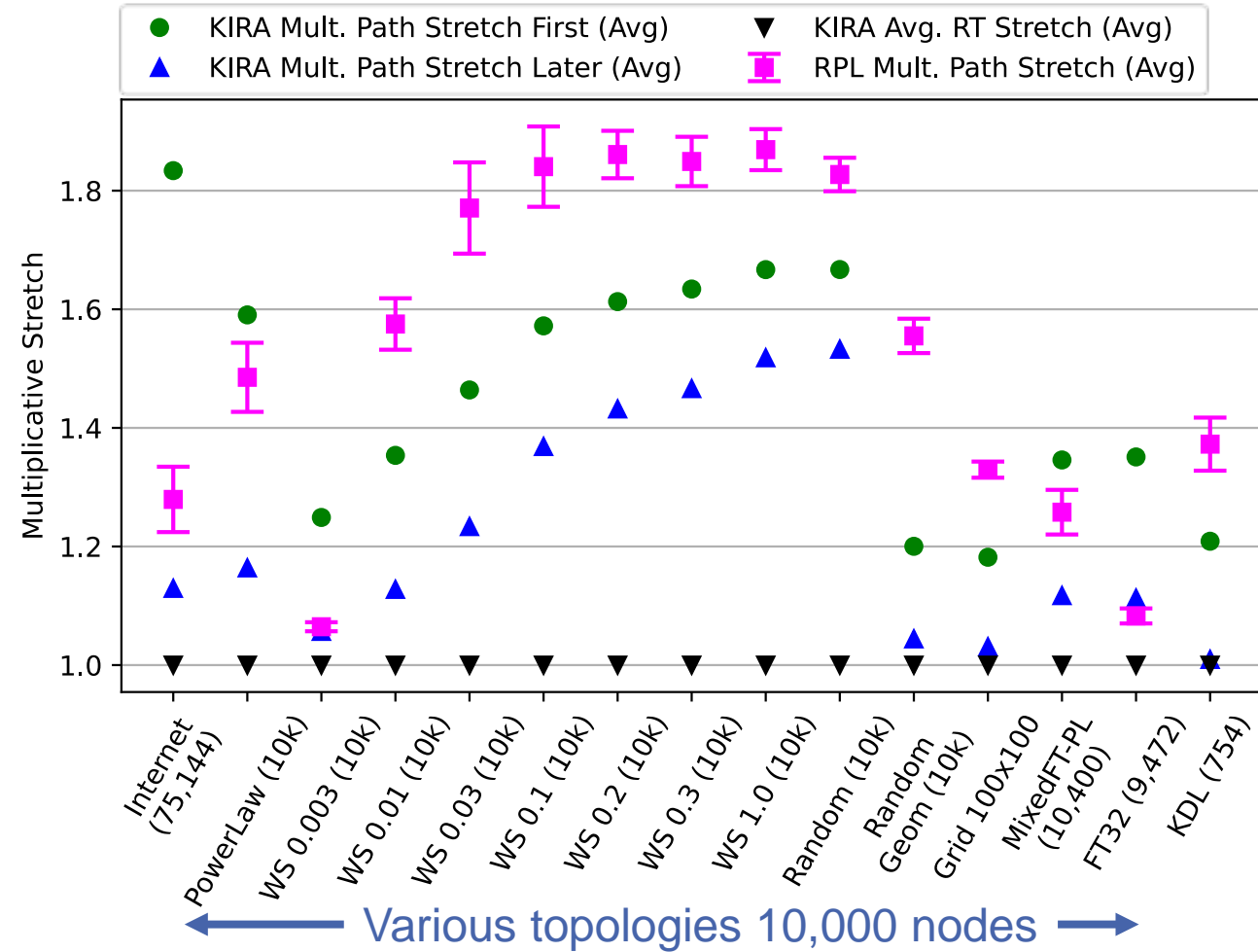
- Fast Failover
 - Fast Next Hop Detour (Forwarding Tier)
 - Fast Vicinity Alternatives (Routing Tier)
- Multi-path capable (due to small routing tables)
- Supports different routing metrics
- Built-in route flapping prevention
- Special end-system mode (non-routing)
- Supports Domains (topological, organizational)
- Can provide additional services like
 - Key-Value store, name registration and resolution service (DHT)
 - KeLLy – efficient topology discovery

Architecture



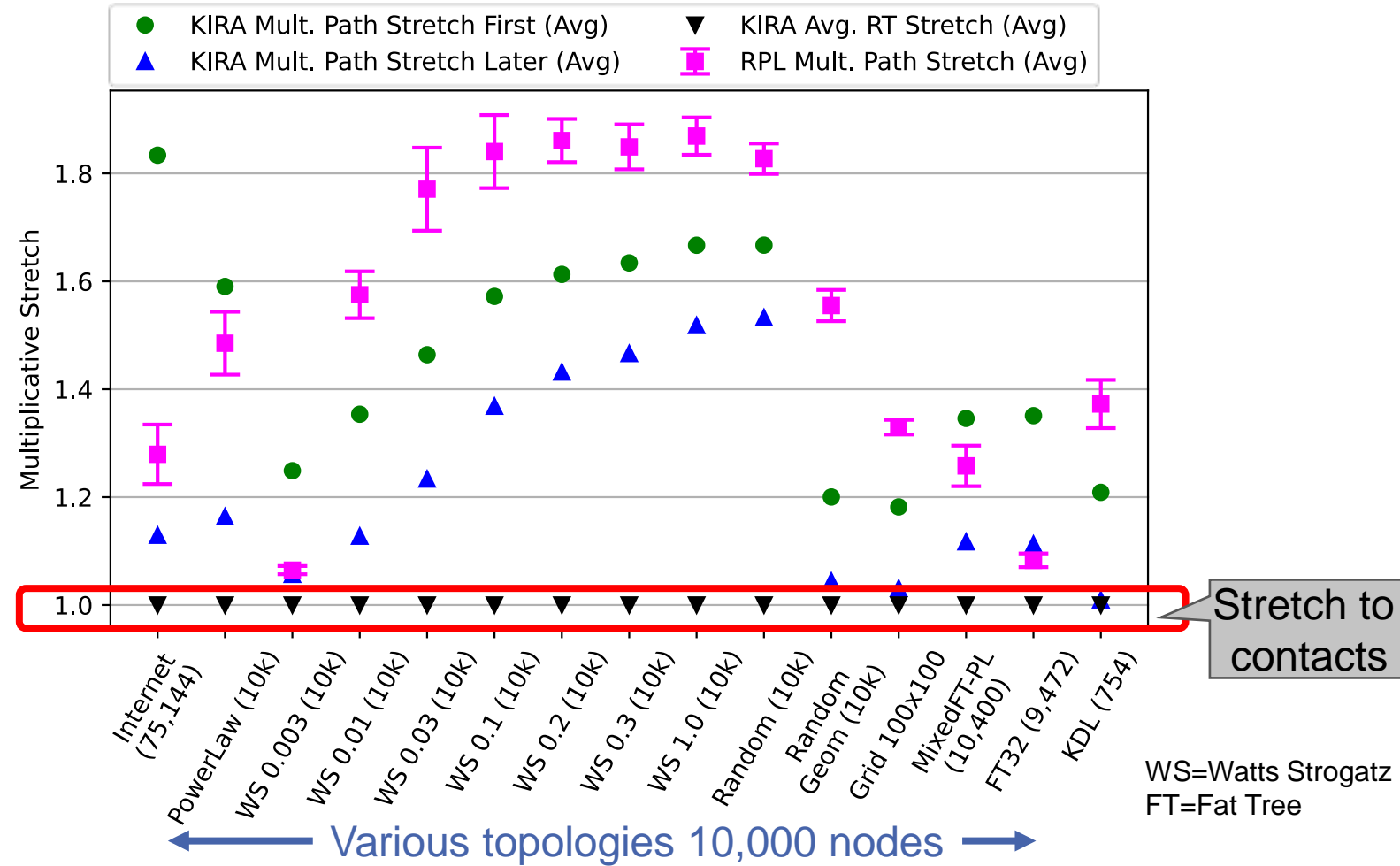
Stretch in Different Fixed Network Topologies

- Multiplicative Stretch
- RPL-ACP:
 - Storing-mode
 - Single DODAG
 - Single DODAG version



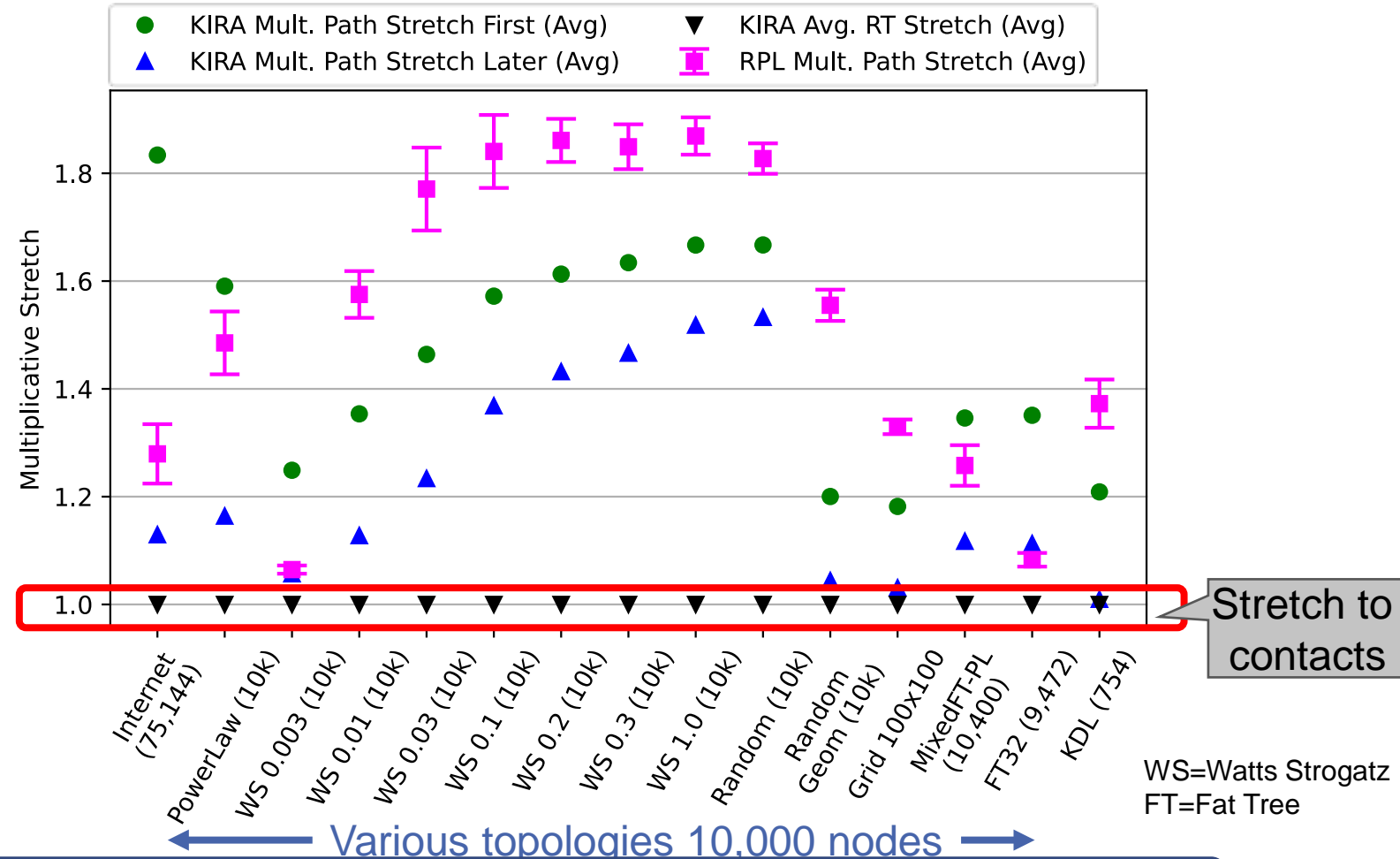
Stretch in Different Fixed Network Topologies

- Multiplicative Stretch
- RPL-ACP:
 - Storing-mode
 - Single DODAG
 - Single DODAG version



Stretch in Different Fixed Network Topologies

- Multiplicative Stretch
- RPL-ACP:
 - Storing-mode
 - Single DODAG
 - Single DODAG version



Low stretch across various topologies + Shortest paths to contacts

The Plan...

■ Status Quo

- Networks and their operations will become more complex
- Increasing number of interdependencies with growing scale (RFC 3439)
- Root cause analysis extremely difficult
- Need resilient operations

■ Let's

- make the Internet better...
- make KIRA an IETF standard so that it can be built into networked devices
- have a resilient connectivity for network control that
 - you can rely on (cannot be broken by misconfiguration)
 - you can always use to restore your network to well-known states

Please Support KIRA!

- Internet-Draft <https://datatracker.ietf.org/doc/draft-bless-rtgwg-kira/>
 - Updated to -03: added also end-system mode
 - Please provide feedback!
 - Want IETF expertise
- Running Code available
 - Hackathon IETF 123
 - Native Routing Daemon Linux (Rust)
 - Zero-touch IPv6 Connectivity
 - Forwarding Tier uses nftables
 - Alternative eBPF implementation ongoing



<https://s.kit.edu/KIRA>

Q&A

Next Steps

Next IETF Steps for KIRA

- Raise community interest: spread the word...
- Move Base Internet-Draft forward (at least to experimental) in RTGWG
- Need feedback!
- Need implementers!
- Think big: WG Topics
 - Base Specification
 - different encapsulation possibilities, Domain-ID in forwarding tier, replacement of hash function
 - Security mechanisms
 - „API“ + protocols for supporting services (e.g., DHT, KeLLy, etc.)
 - Specially tweaked versions for IoT, MANETs etc.

Backup Slides

KIRA – Main Components

■ Routing Tier → connectivity

R²/Kad

- ID-based addresses
- Source routing
- On top of link layer

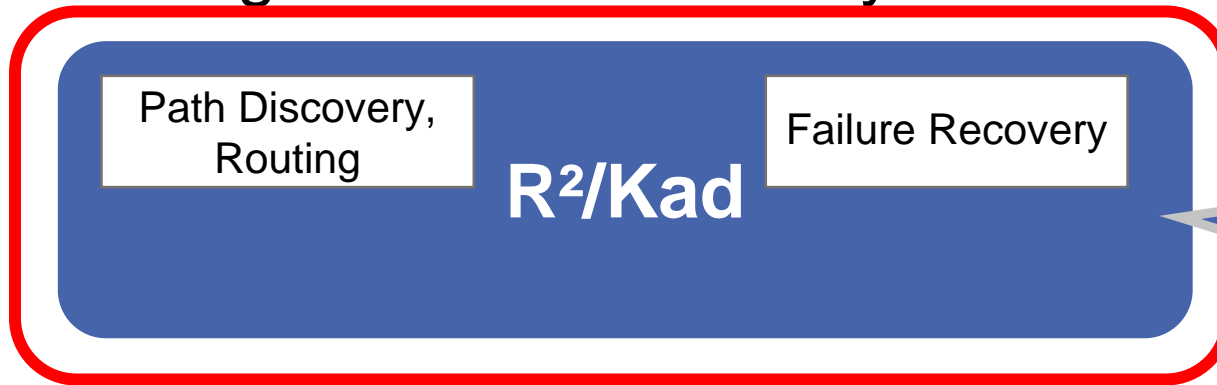
■ Forwarding Tier → optimization

PathID-based Forwarding

- Eliminates source routing
- Label switching approach
- Reduces overhead

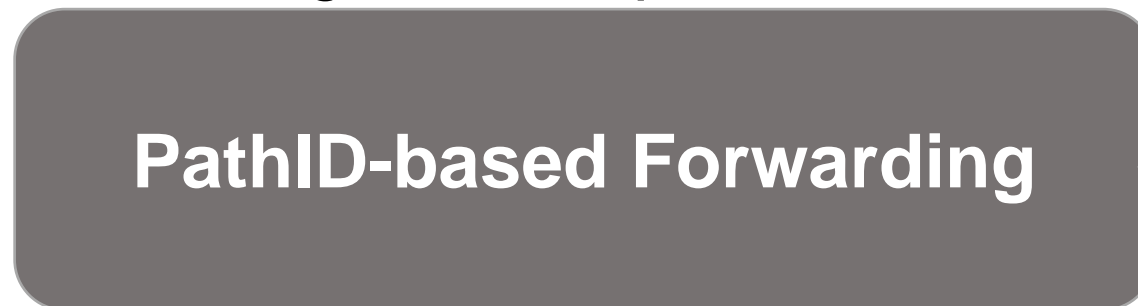
KIRA – Main Components

■ Routing Tier → connectivity



- ID-based addresses
- Source routing
- On top of link layer

■ Forwarding Tier → optimization



- Eliminates source routing
- Label switching approach
- Reduces overhead

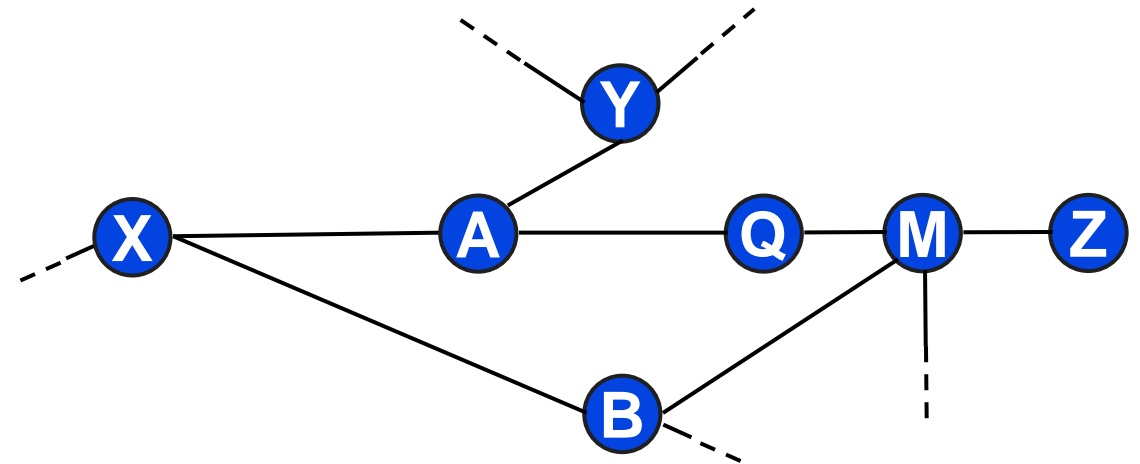
R²/Kad – Path Discovery

- Each node
 - randomly chooses its **NodeID** (Overlay)

Path Discovery,
Routing

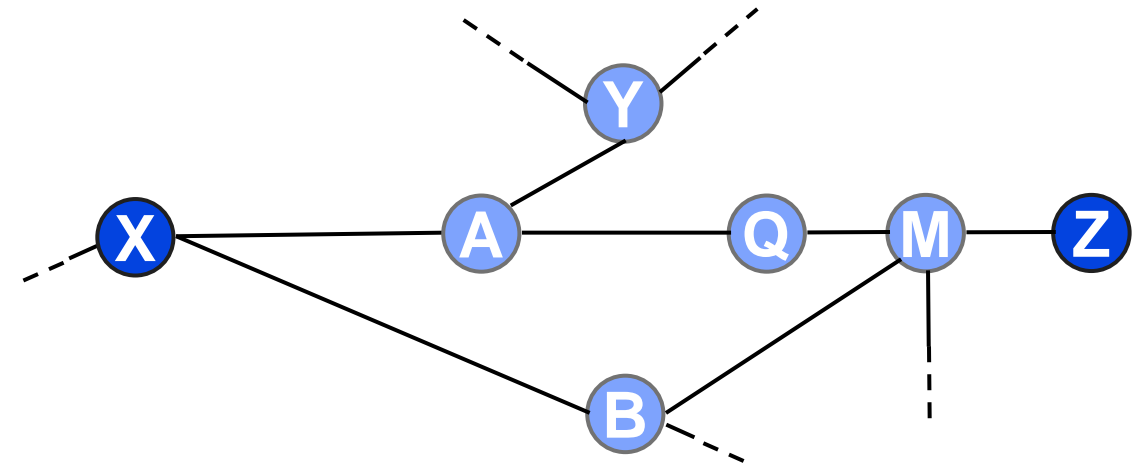
R²/Kad

Failure
Recovery



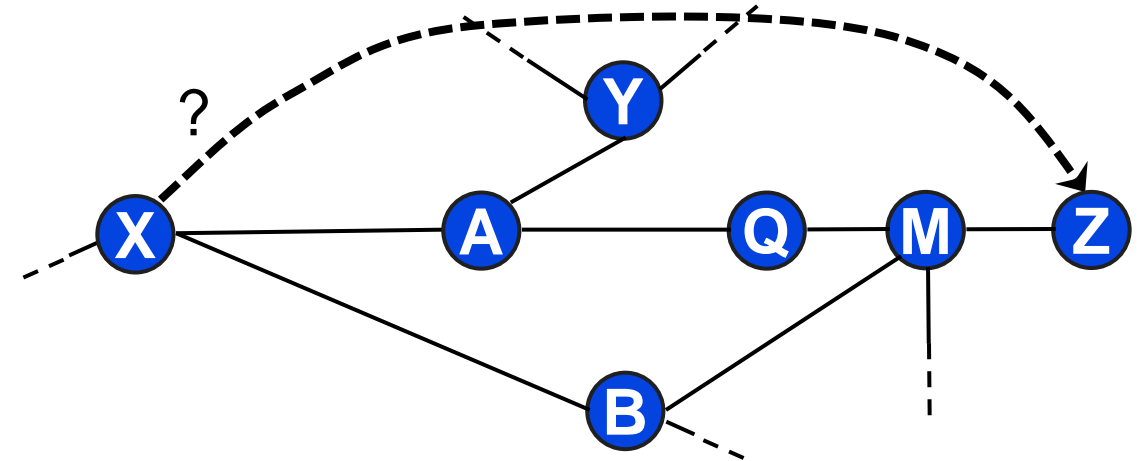
R²/Kad – Path Discovery

- Each node
 - randomly chooses its **NodeID** (Overlay)
 - explores its **2-hop vicinity** (Underlay)
 - X learns contacts A, Y, Q, B, M, ...



R²/Kad – Path Discovery

- Each node
 - randomly chooses its **NodeID** (Overlay)
 - explores its **2-hop vicinity** (Underlay)
 - X learns contacts A, Y, Q, B, M, ...
- X: path to Z?



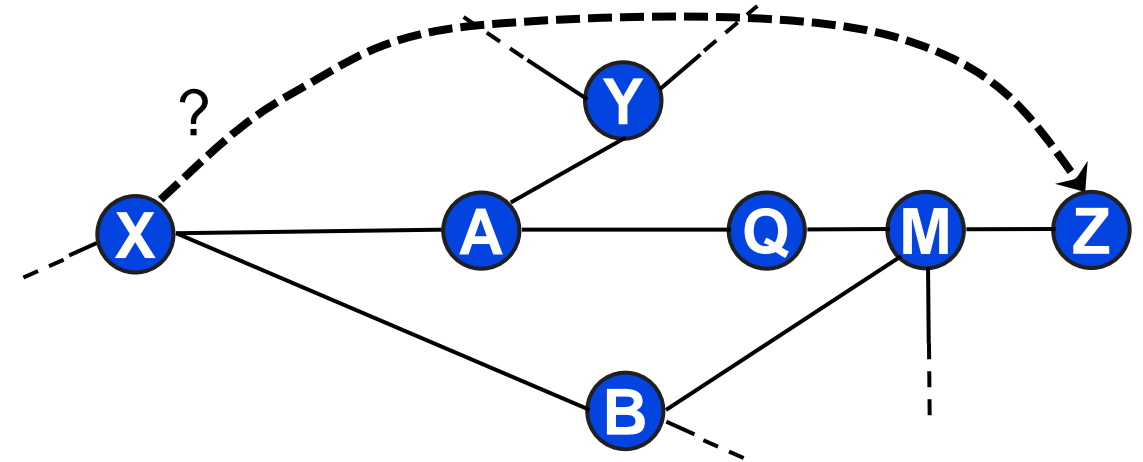
R²/Kad – Path Discovery

Path Discovery,
Routing

R²/Kad

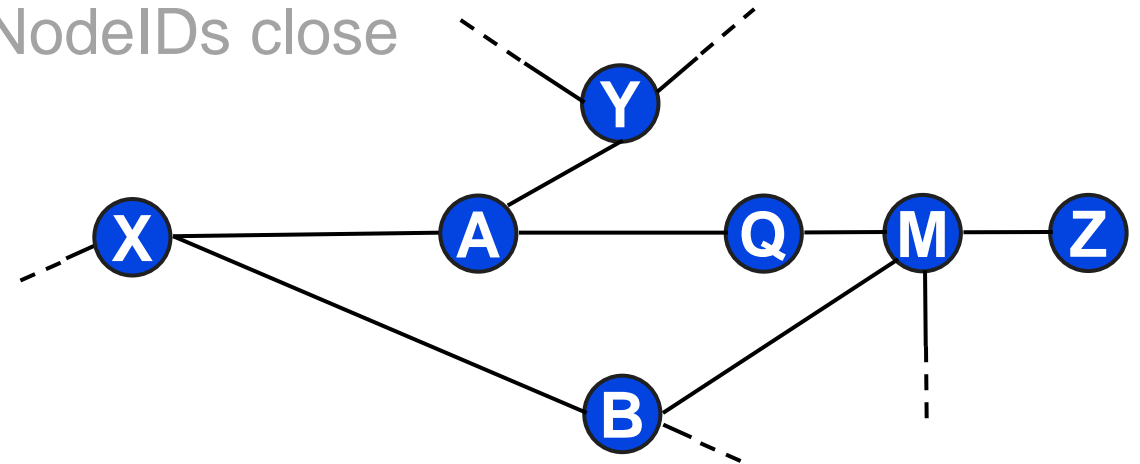
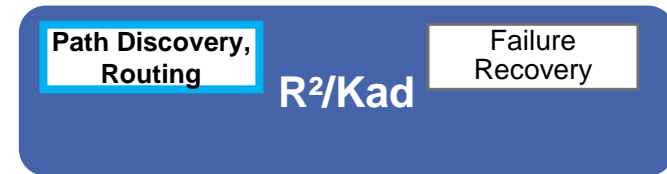
Failure
Recovery

- Each node
 - randomly chooses its **NodeID** (Overlay)
 - explores its **2-hop vicinity** (Underlay)
 - X learns contacts A, Y, Q, B, M, ...
- X: path to Z?
- Approach:
 - construct underlay routes
 - by using the **NodeID-based overlay**
 - Source route to contact that is **ID-wise closest** to destination NodeID (→ recursively)
 - Distance of NodeIDs: **XOR metric** $d(X, Y) = X \oplus Y$
 - Longer shared prefix → closer



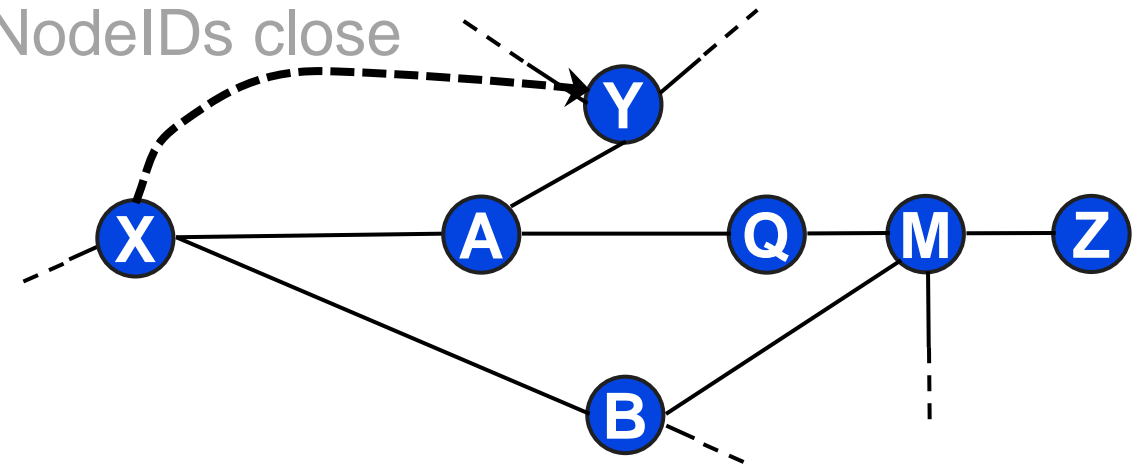
R²/Kad – Path Discovery Example

- X sends FindNodeReq to contact **closest** to NodeID Z
 - Example: letters close in alphabet \leftrightarrow NodeIDs close



R²/Kad – Path Discovery Example

- X sends FindNodeReq to contact **closest** to NodeID Z
 - Example: letters close in alphabet \leftrightarrow NodeIDs close
 - Next (overlay) hop: Y



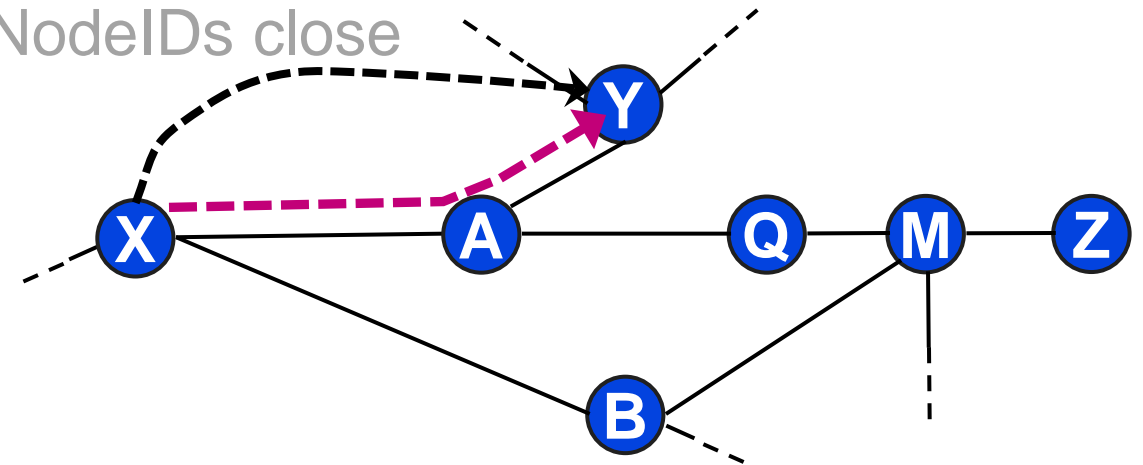
Path Discovery,
Routing

R²/Kad

Failure
Recovery

R²/Kad – Path Discovery Example

- X sends FindNodeReq to contact **closest** to NodeID Z
 - Example: letters close in alphabet \leftrightarrow NodeIDs close
 - Next (overlay) hop: Y
- X \rightarrow Y via **source route** <A>

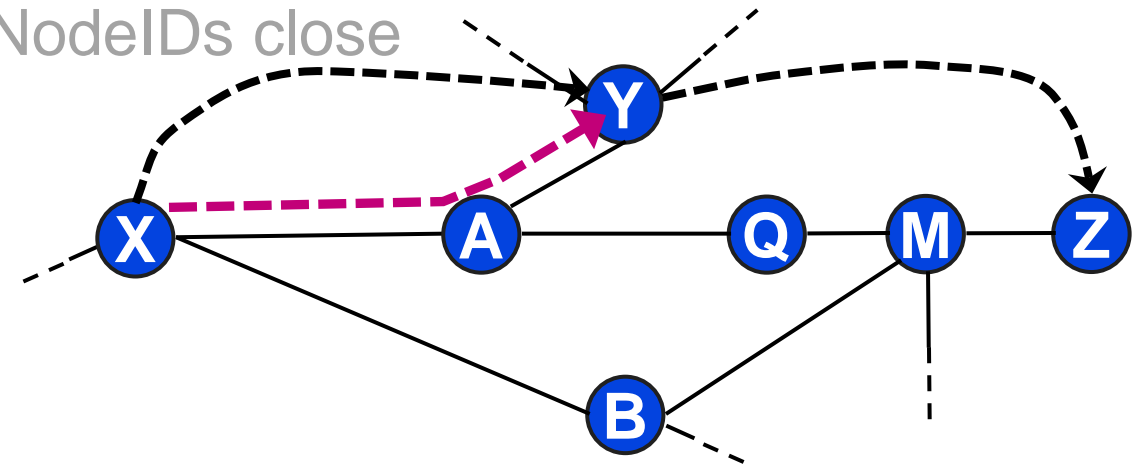

Path Discovery,
Routing

R²/Kad

Failure
Recovery

R²/Kad – Path Discovery Example

- X sends FindNodeReq to contact **closest** to NodeID Z
 - Example: letters close in alphabet ↔ NodeIDs close
 - Next (overlay) hop: Y
- X → Y via **source route** <A>
- Assume Y knows Z already


Path Discovery,
Routing

R²/Kad

Failure
Recovery

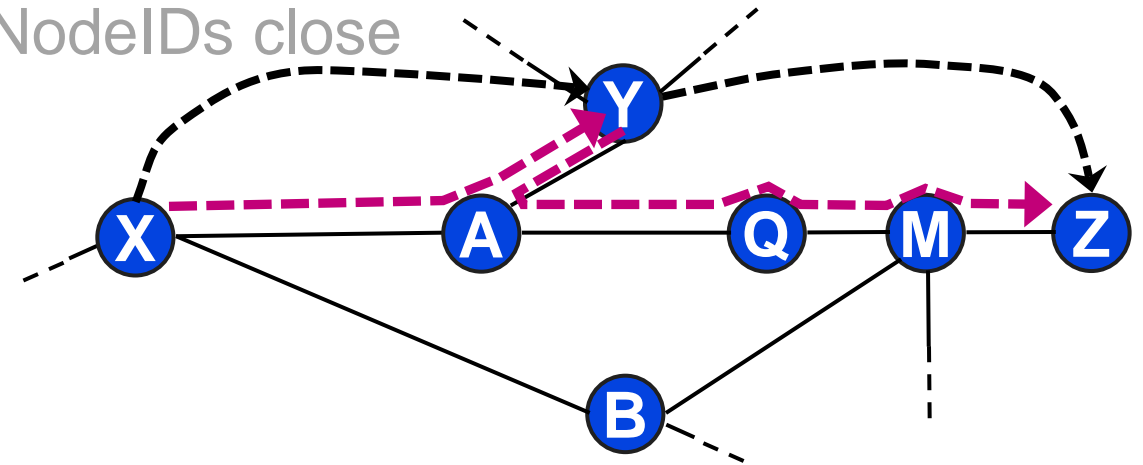
R²/Kad – Path Discovery Example

Path Discovery,
Routing

R²/Kad

Failure
Recovery

- X sends FindNodeReq to contact **closest** to NodeID Z
 - Example: letters close in alphabet ↔ NodeIDs close
 - Next (overlay) hop: Y
- X → Y via **source route** <A>
- Assume Y knows Z already
- Y → Z via **source route** <A,Q,M>
- FindNodeReq records complete route <X,A,Y,A,Q,M>
 - incurs path stretch: $\frac{|\text{selected path}|}{|\text{shortest path}|}$

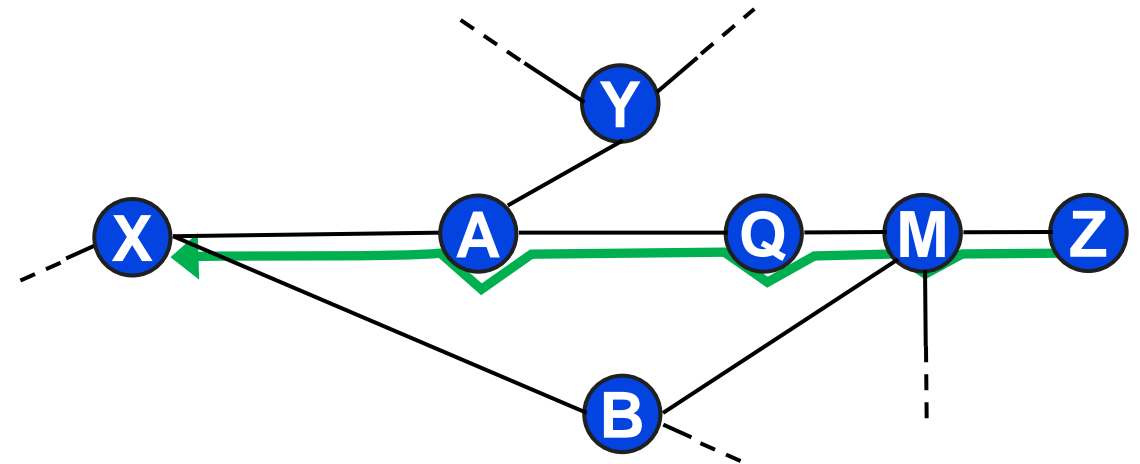


R²/Kad – Path Discovery Example

- Shortened recorded route <A,Q,M> is returned to X in FindNodeRsp

Path Discovery,
Routing

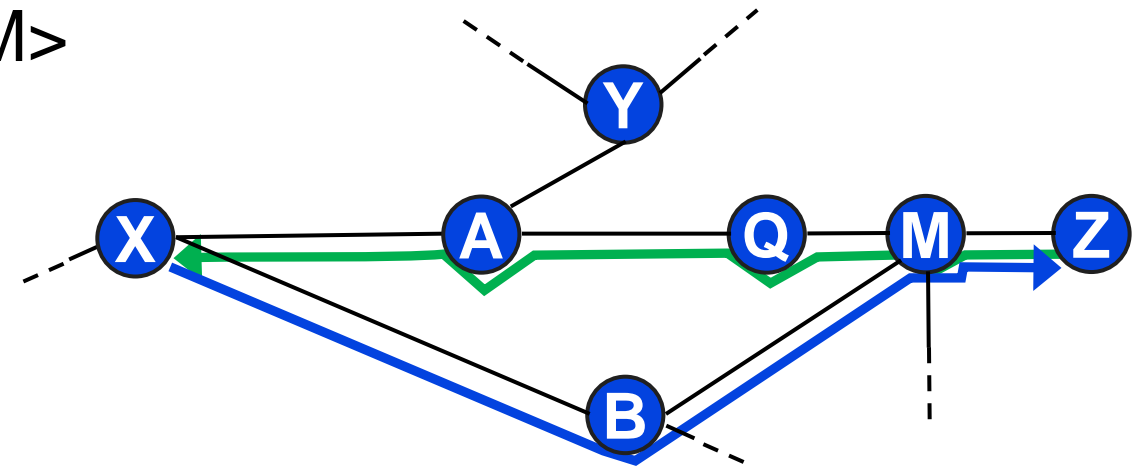
R²/Kad

Failure
Recovery


R²/Kad – Path Discovery Example

- Shortened recorded route $\langle A, Q, M \rangle$ is returned to X in FindNodeRsp
- Later packets use shorter route $\langle B, M \rangle$
 - if X already knows M via $\langle B \rangle$

Initial stretch can be reduced for later packets!



Path Discovery,
Routing

R²/Kad

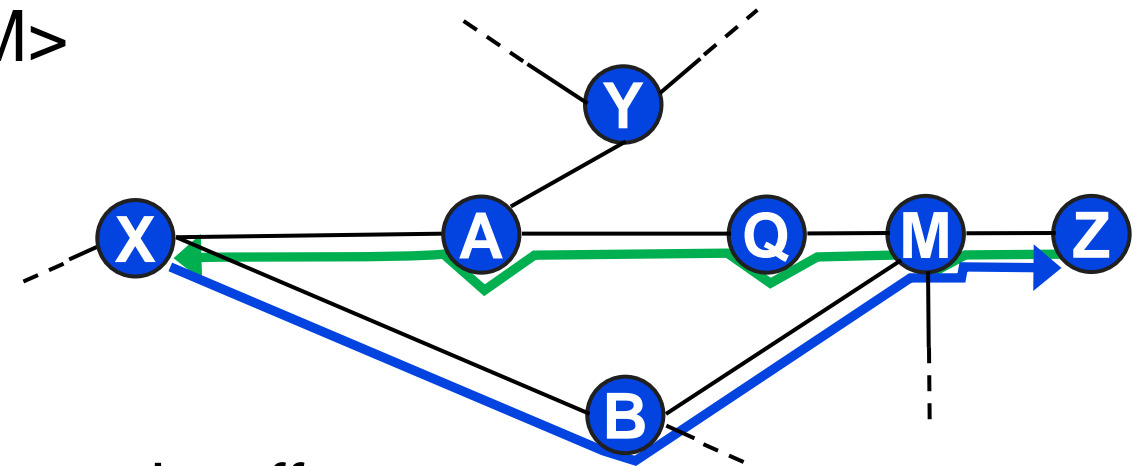
Failure
Recovery

R²/Kad – Path Discovery Example

- Shortened recorded route $\langle A, Q, M \rangle$ is returned to X in FindNodeRsp
- Later packets use shorter route $\langle B, M \rangle$
 - if X already knows M via $\langle B \rangle$

Initial stretch can be reduced for later packets!

- R²/Kad offers flexible memory/stretch trade-off...



Path Discovery,
Routing

R²/Kad

Failure
Recovery

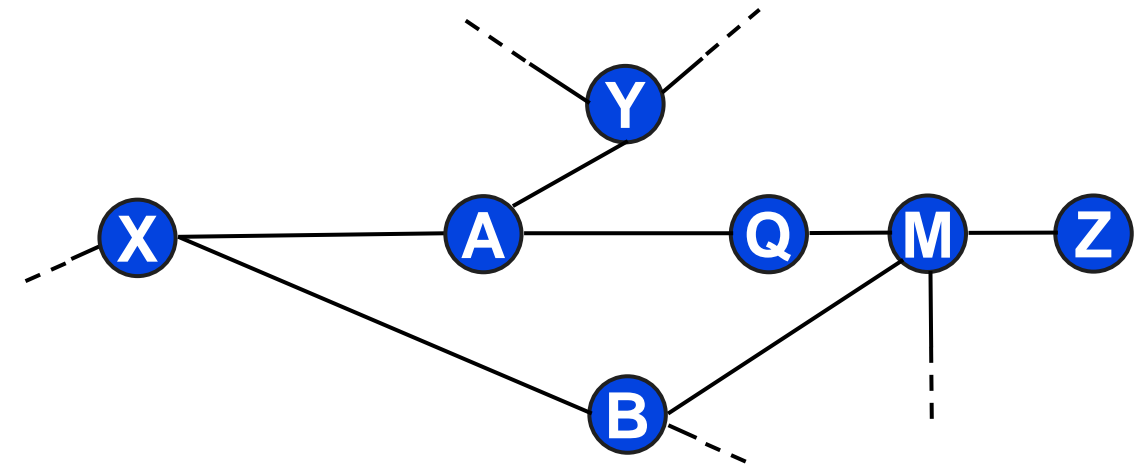
R²/Kad – Dynamics: Rediscovery Procedure

Path Discovery,
Routing

R²/Kad

Failure
Recovery

- Detection of node/link failure in the underlay



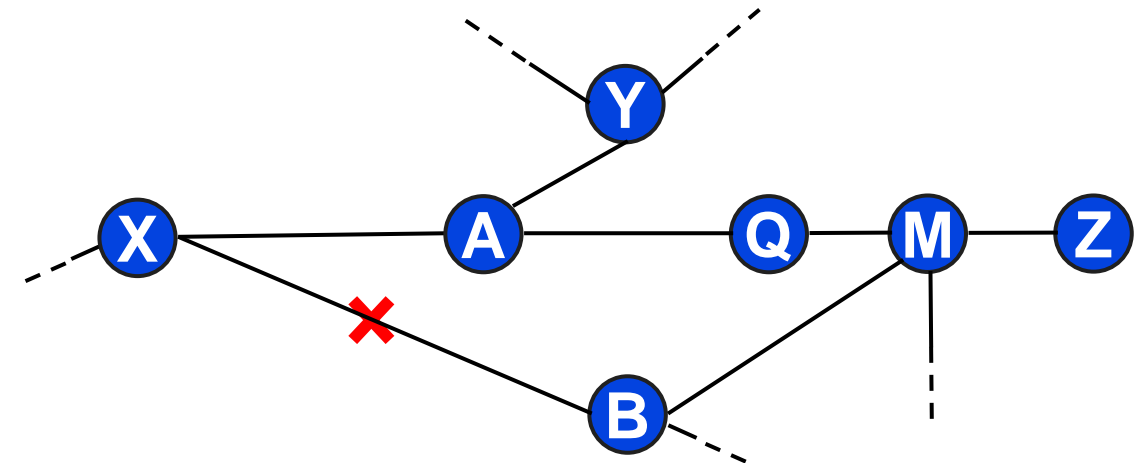
R²/Kad – Dynamics: Rediscovery Procedure

Path Discovery,
Routing

R²/Kad

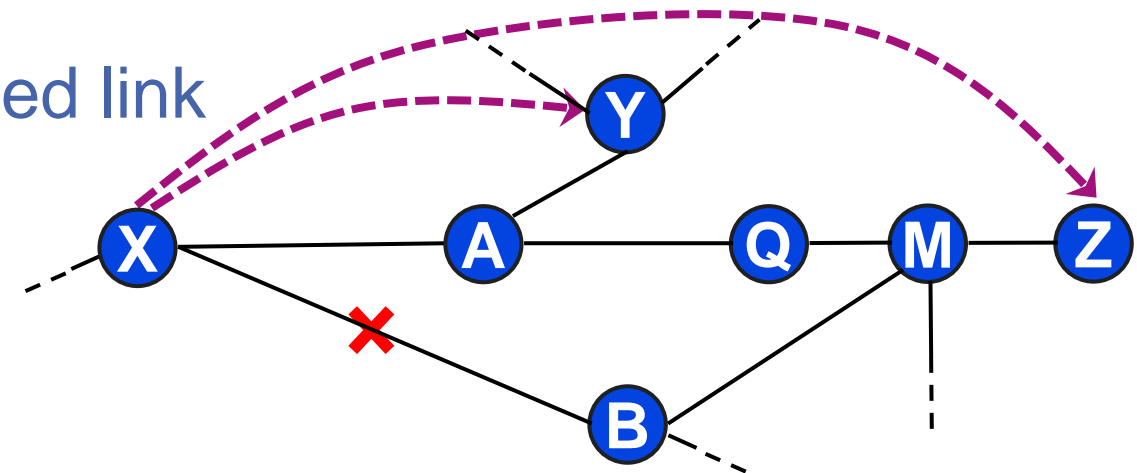
Failure
Recovery

- Detection of node/link failure in the underlay



R²/Kad – Dynamics: Rediscovery Procedure

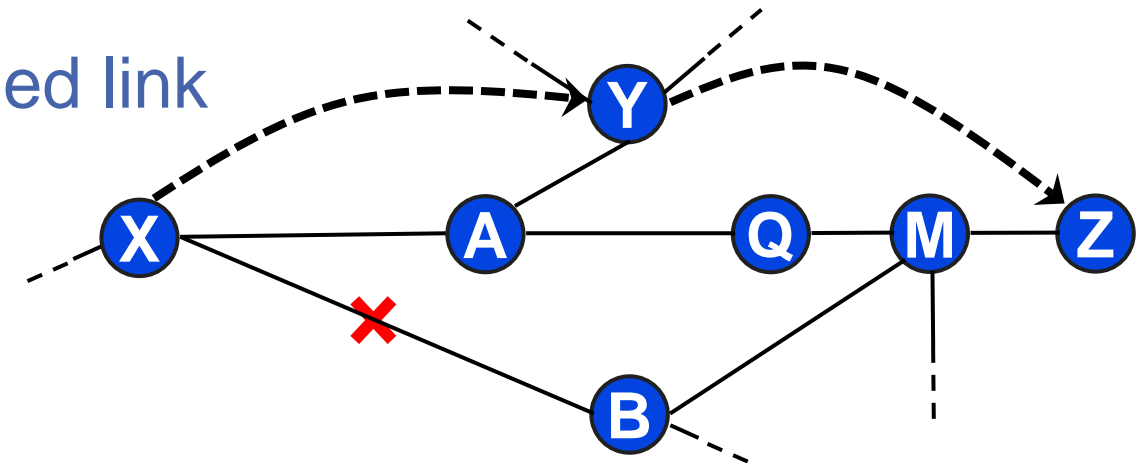
- Detection of node/link failure in the underlay
- Two step strategy
 - 1.) **inform** ID-wise neighbors about **failed link**
 - 2.) ...



R²/Kad – Dynamics: Rediscovery Procedure



- Detection of node/link failure in the underlay
- Two step strategy
 - 1.) **inform** ID-wise neighbors about **failed link**
 - 2.) **rediscover** alternative paths via overlay routes (includes “**Not Via**” information)
- Validity
 - **State sequence numbers**
 - Path information **age**
- Periodically
 - probe contacts for broken paths
 - lookup own NodeID



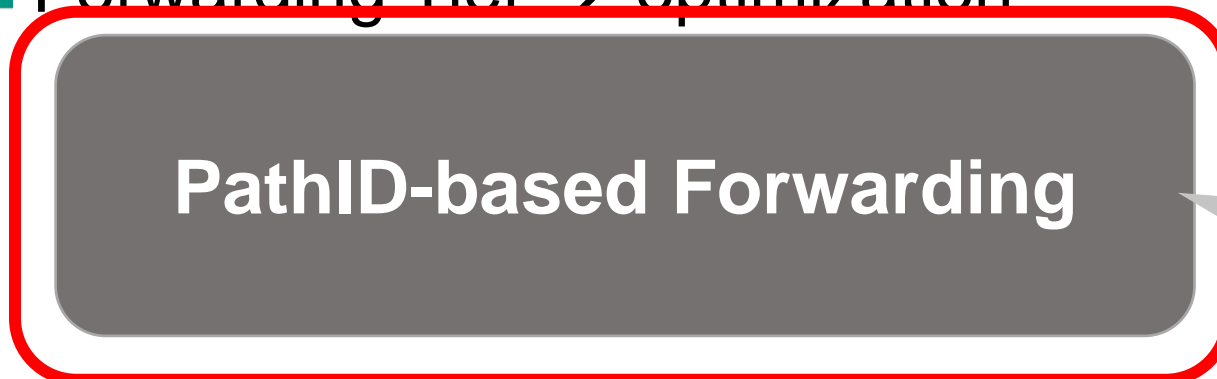
KIRA – Main Components

■ Routing Tier → connectivity



- ID-based addresses
- Source routing
- On top of link layer

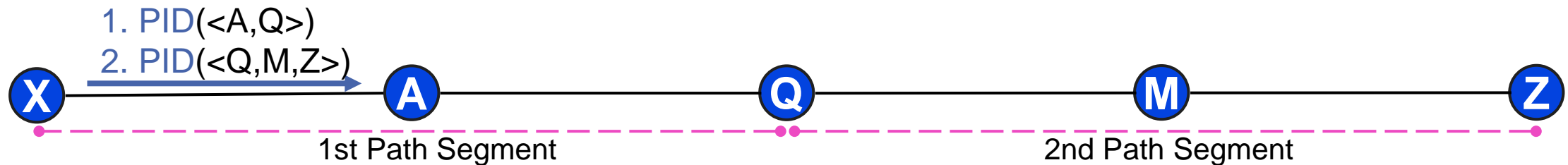
■ Forwarding Tier → optimization



- Label Switching Approach
- Eliminates Source Routing
- Reduces Overhead

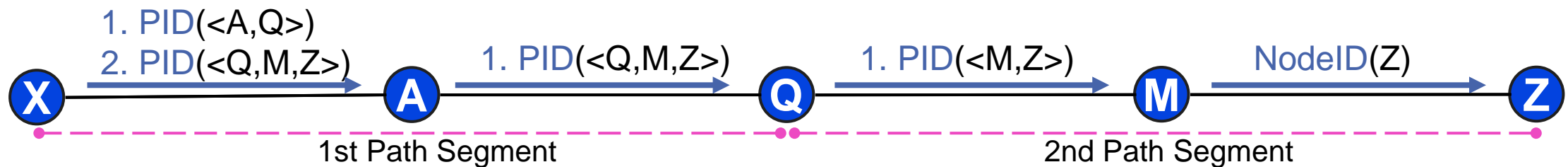
Forwarding Tier – Fast Forwarding

- Get rid of source routes for control plane traffic
 - Reduce per packet overhead
- Approach: **replace source routes with PathIDs**
 - $\text{PathID}(\langle A, Q, M, Z \rangle) = \text{Hash}(A \mid Q \mid M \mid Z)$
- Use PathID as unique label for path segment → Label Switching



Forwarding Tier – Fast Forwarding

- Get rid of source routes for control plane traffic
 - Reduce per packet overhead
- Approach: **replace source routes with PathIDs**
 - $\text{PathID}(\langle A, Q, M, Z \rangle) = \text{Hash}(A \mid Q \mid M \mid Z)$
- Use PathID as unique label for path segment → Label Switching



- Precalculate PathIDs for 2-hop (physical) vicinity
- Explicit **path setup** for paths ≥ 6 hops

State Sequence Numbers

■ Per Node: State Sequence Number

- reflects changes in node's physical neighbor set
- Link down
- Link up (also detecting a new node)

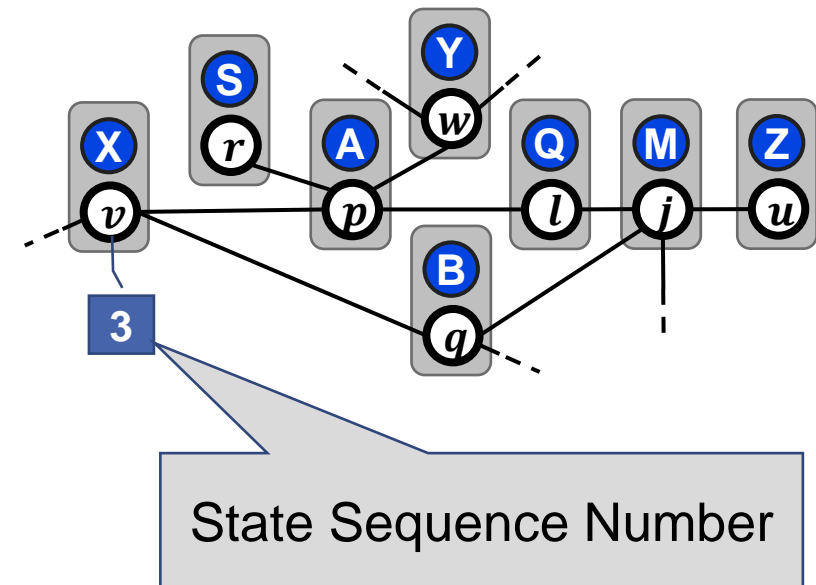
■ 32-bit

- Wrap around and special comparison:
 - $s < s' \bmod 2^{32}$ if $0 < (s' - s) \bmod 2^{32} < 2^{31}$

■ Get periodically synchronized

■ Node crashes

- Node either uses new NodeID after restart
- or, node stores NodeID and State Sequence Number across restart



State Sequence Numbers

■ Per Node: State Sequence Number

- reflects changes in node's physical neighbor set
- Link down
- Link up (also detecting a new node)

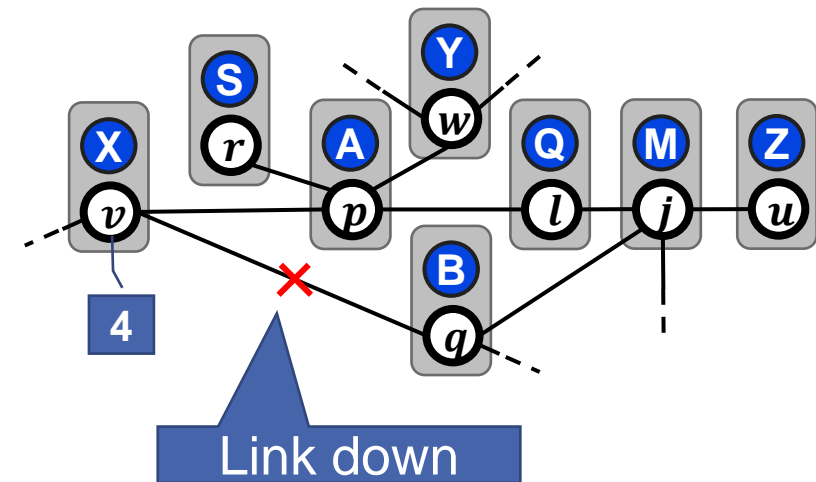
■ 32-bit

- Wrap around and special comparison:
 - $s < s' \bmod 2^{32}$ if $0 < (s' - s) \bmod 2^{32} < 2^{31}$

■ Get periodically synchronized

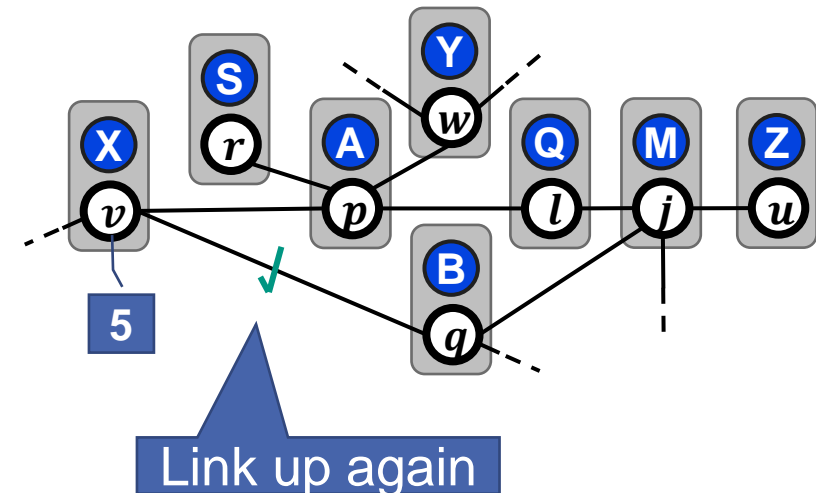
■ Node crashes

- Node either uses new NodeID after restart
- or, node stores NodeID and State Sequence Number across restart



State Sequence Numbers

- Per Node: **State Sequence Number**
 - reflects changes in node's physical neighbor set
 - Link down
 - Link up (also detecting a new node)
- 32-bit
 - Wrap around and special comparison:
 - $s < s' \bmod 2^{32}$ if $0 < (s' - s) \bmod 2^{32} < 2^{31}$
- Get periodically synchronized
- Node crashes
 - Node either uses new NodeID after restart
 - or, node stores NodeID and State Sequence Number across restart



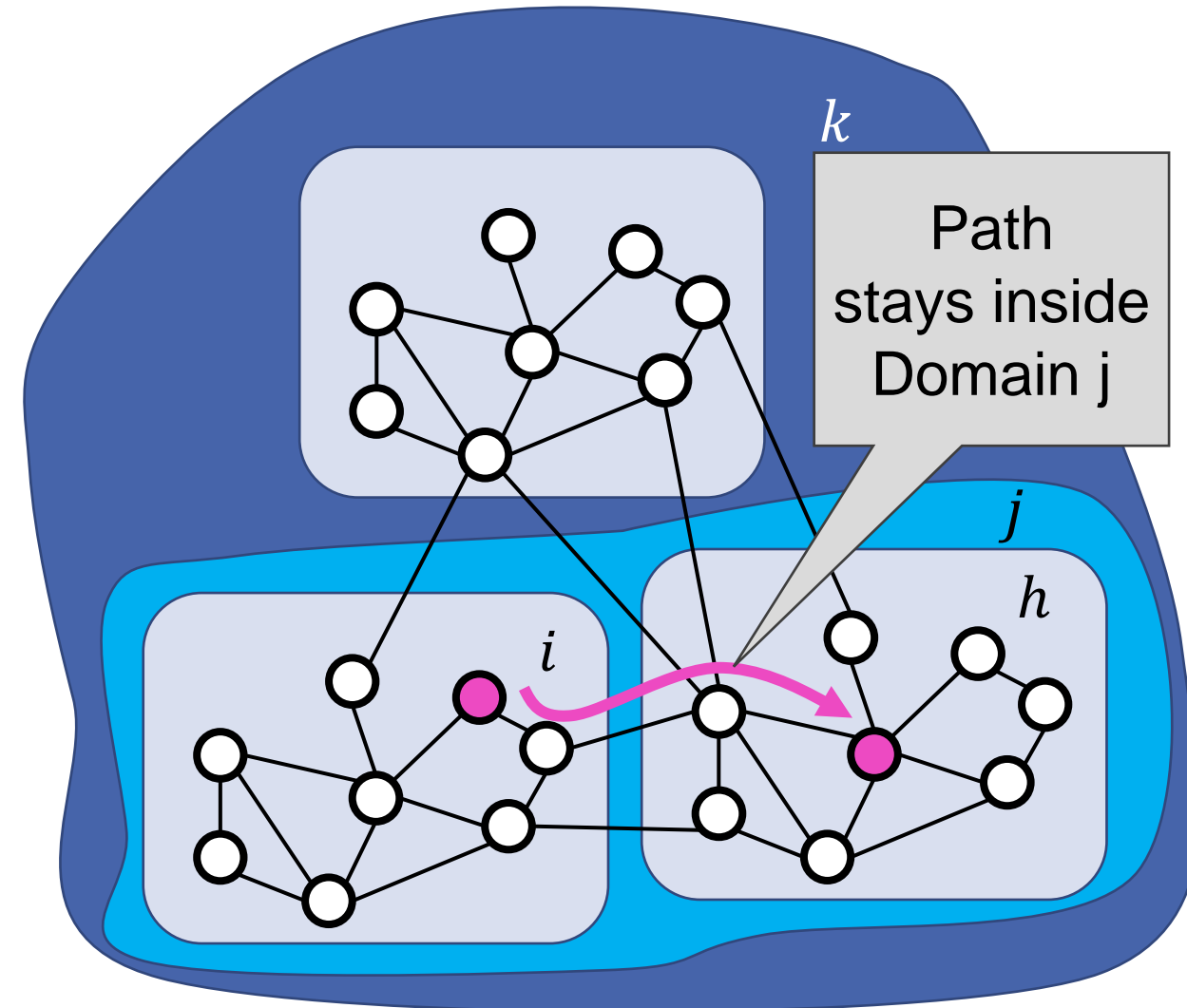
End-system Mode

- End-systems do not route, but may be multi-homed and mobile
- Reduce overhead by not transmitting routing updates to/from end-systems
- Routers are responsible to keep information on end-system reachability

KIRA – Domain Scopes

■ Domain Scopes

- Global, Organizational (e.g., ASes), Topological
- KIRA nodes keep their NodeID across domains!



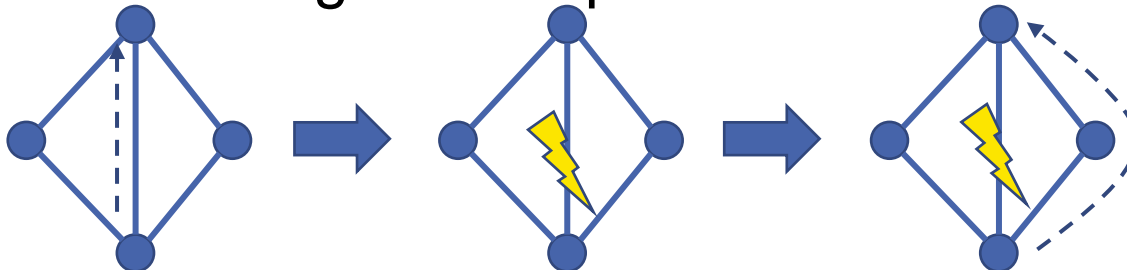
Recent Changes (1)

■ Path Validation:

- (Partial) source routes are considered to be valid
- Routing information from other nodes must be validated before actual use by a ProbeReq/-Rsp (Proposed Path → Active Path)

■ Fast Vicinity Alternatives

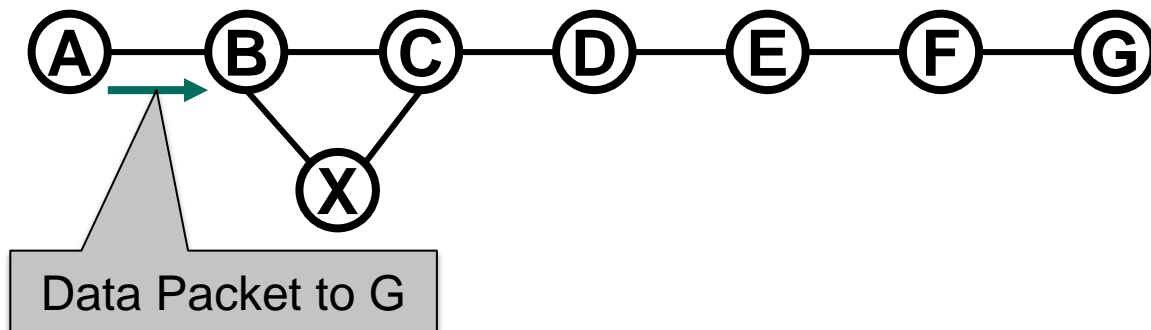
- KIRA maintains a **vicinity graph (full 2-hop underlay vicinity)** to precompute PathIDs
- If underlay neighbor fails, try to find alternative route in vicinity graph
- No convergence required due to source routing (PathIDs)



Recent Changes (2)

■ Fast Next Hop Detour

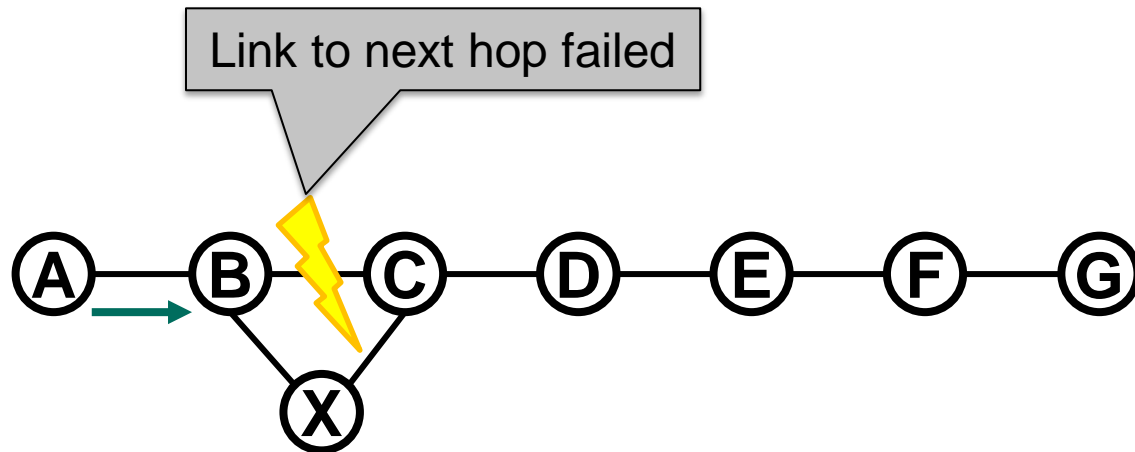
- Precalculate 2-hop backup paths for direct underlay neighbors
- Fast reaction in the **Forwarding Tier** if next hop fails
- Adds one more PathID for detour, gets removed at next hop



Recent Changes (2)

■ Fast Next Hop Detour

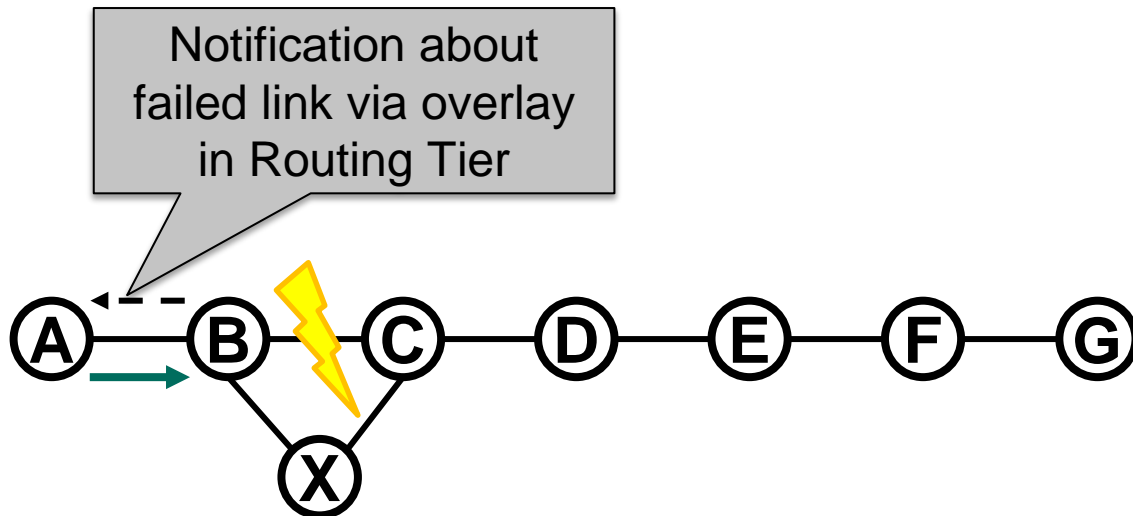
- Precalculate 2-hop backup paths for direct underlay neighbors
- Fast reaction in the **Forwarding Tier** if next hop fails
- Adds one more PathID for detour, gets removed at next hop



Recent Changes (2)

■ Fast Next Hop Detour

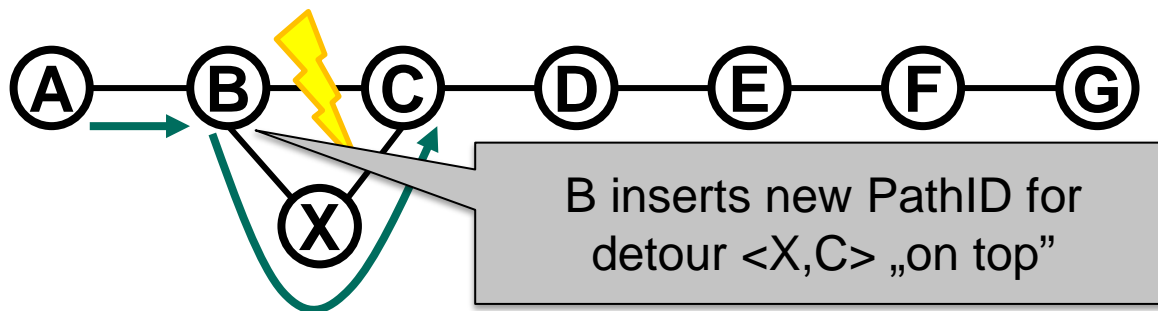
- Precalculate 2-hop backup paths for direct underlay neighbors
- Fast reaction in the **Forwarding Tier** if next hop fails
- Adds one more PathID for detour, gets removed at next hop



Recent Changes (2)

■ Fast Next Hop Detour

- Precalculate 2-hop backup paths for direct underlay neighbors
- Fast reaction in the **Forwarding Tier** if next hop fails
- Adds one more PathID for detour, gets removed at next hop



- No recomputation of outgoing PathID at B necessary
- No change in general forwarding behavior necessary
- No forwarding loops possible, because of continuation on original segment
- Repeated application possible, e.g., if link <X,C> has also a failure