

# *FRANCIS*

## Fast Reaction Algorithms for Network Coordination In Switches

**Presenter: Wenchen Han (UCL)**

Joint work with: Vic Feng, Gregory Schwartzman, Yuliang Li,  
Michael Mitzenmacher, Minlan Yu, **Ran Ben Basat**



HARVARD  
UNIVERSITY

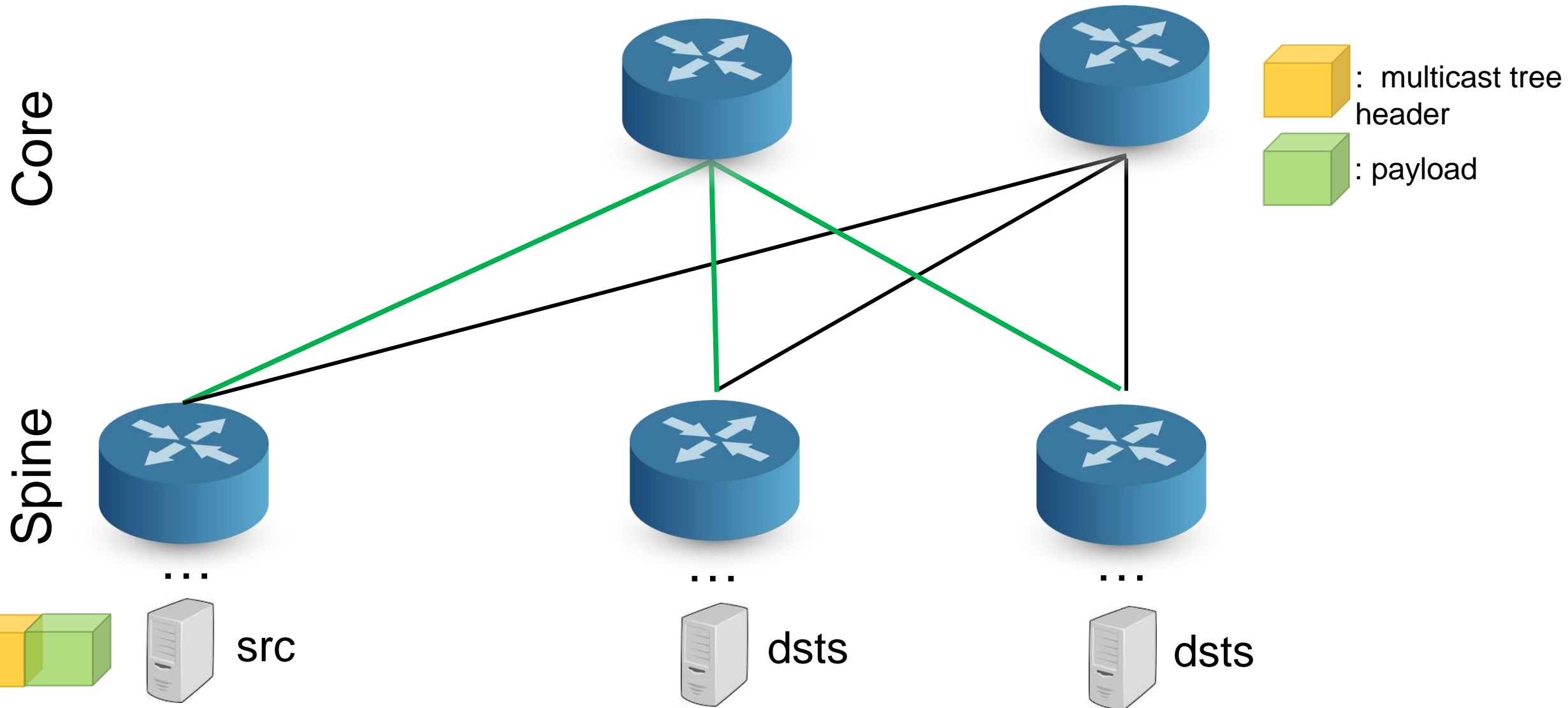


# The need of real-time network changes handling

- **What** are network changes: link/switch failures, congestions, short traffic bursts, etc.
- **Why** real-time: minimizing their impact to applications (multicast, clock-sync, unicast).
- **SDN approaches**:  $O(100ms)$ , inherently too slow for many applications.

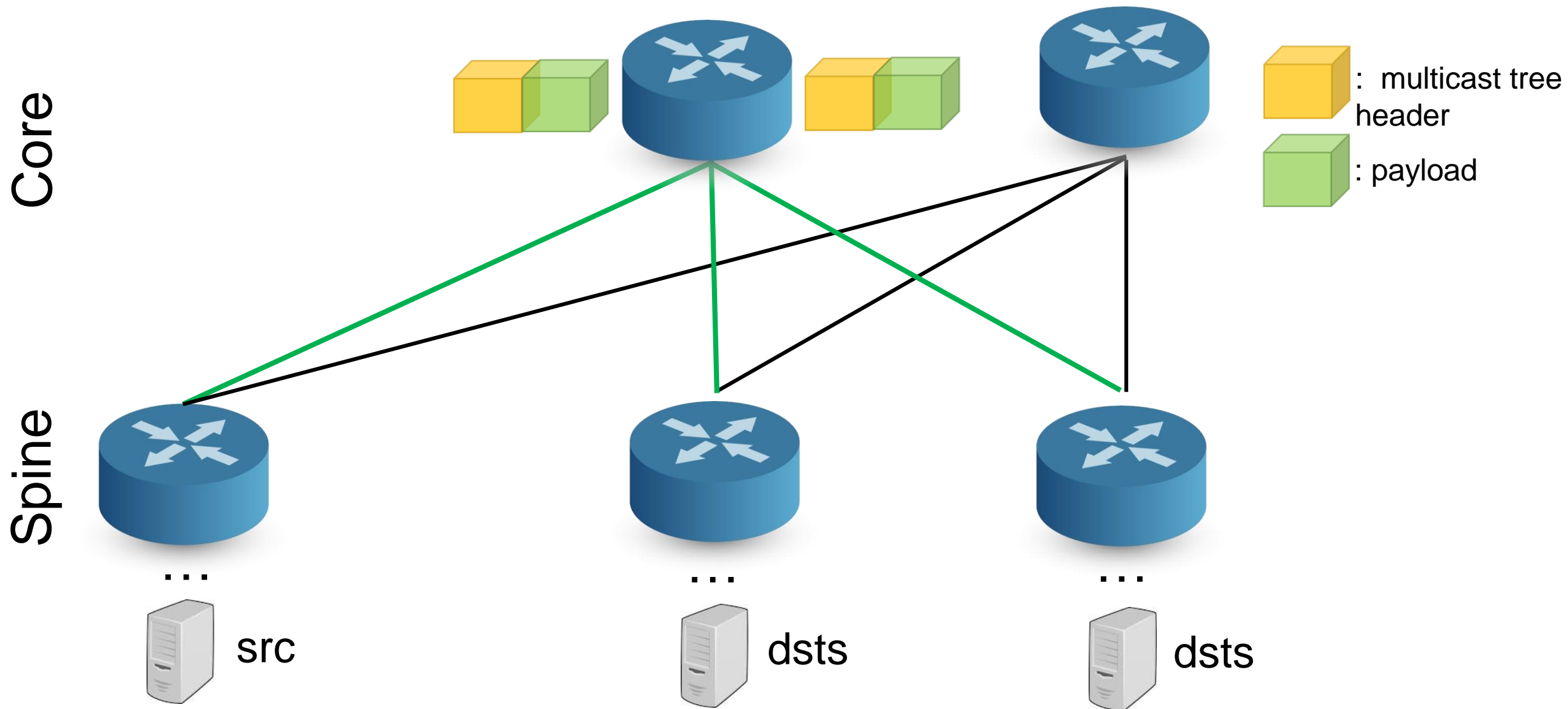


# For example: Elmo [1]'s source-routed multicast



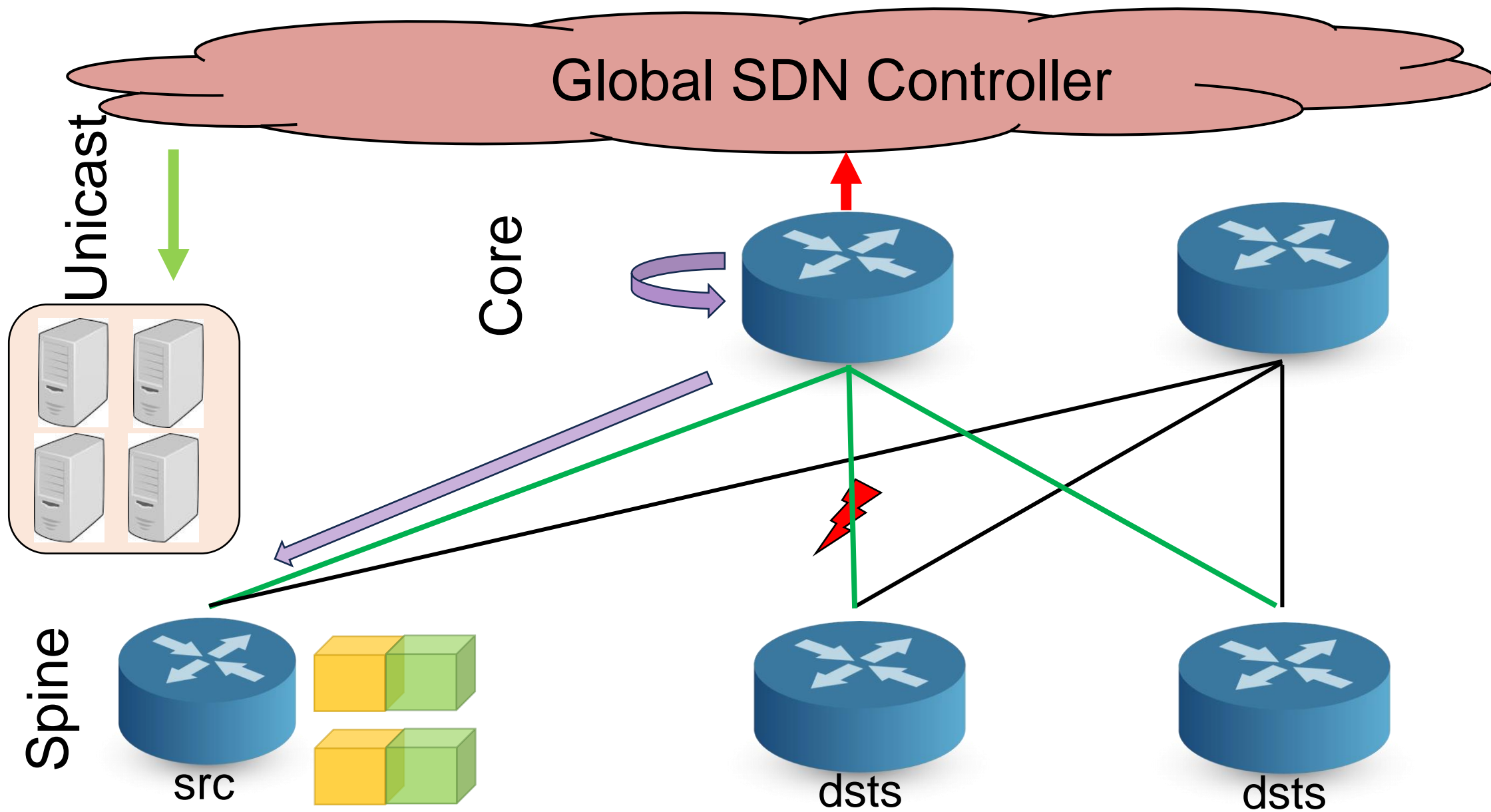
[1] Shahbaz et al. Elmo: Source Routed Multicast for Public Clouds. In SIGCOMM 2019.

# For example: Elmo [1]'s source-routed multicast




# Elmo's SDN-based failure recovery

⌚ 3ms (DP based)



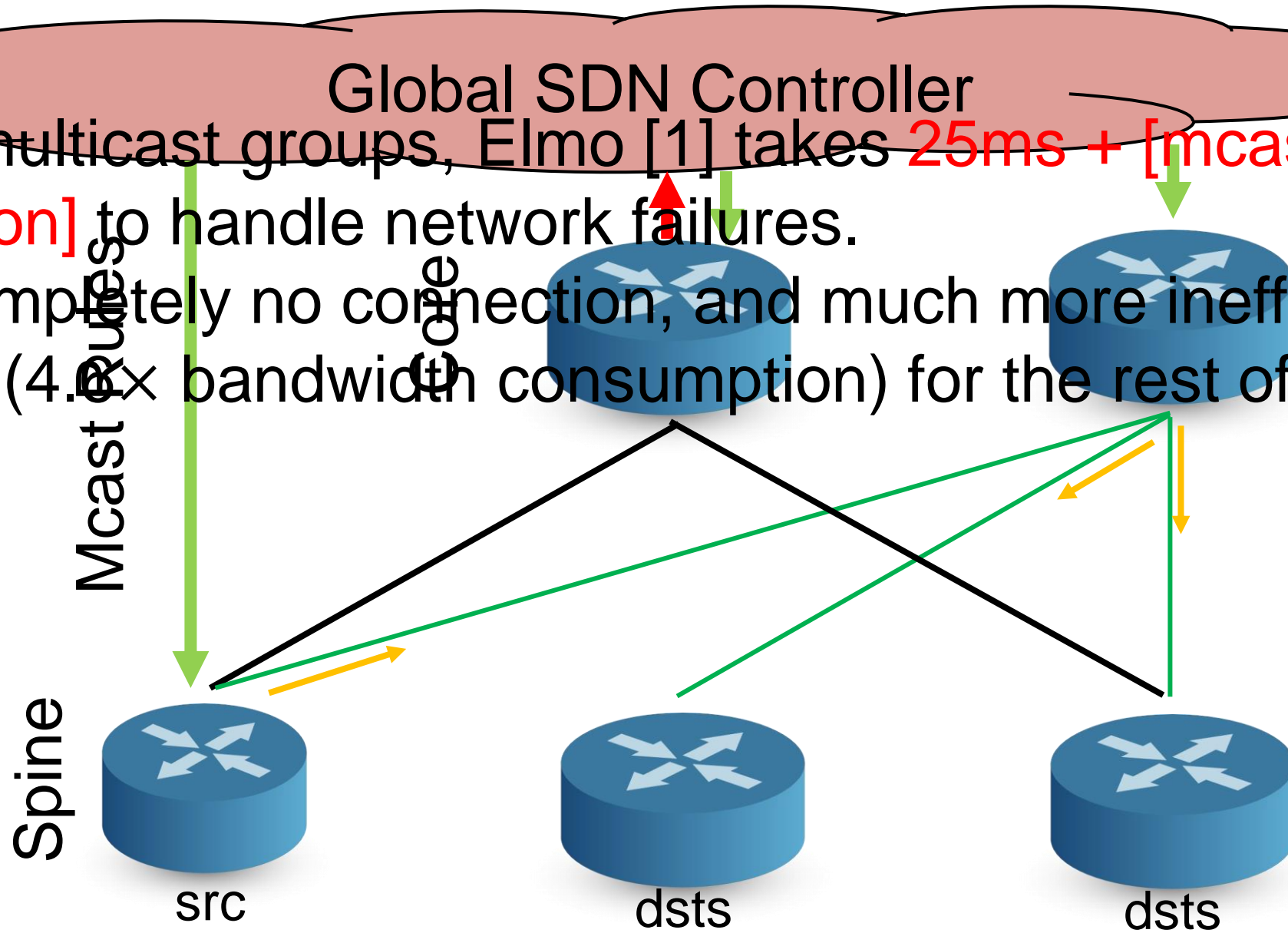
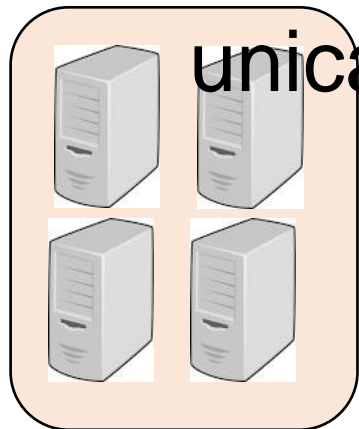
# Elmo's SDN-based failure recovery

  $O(25ms)$

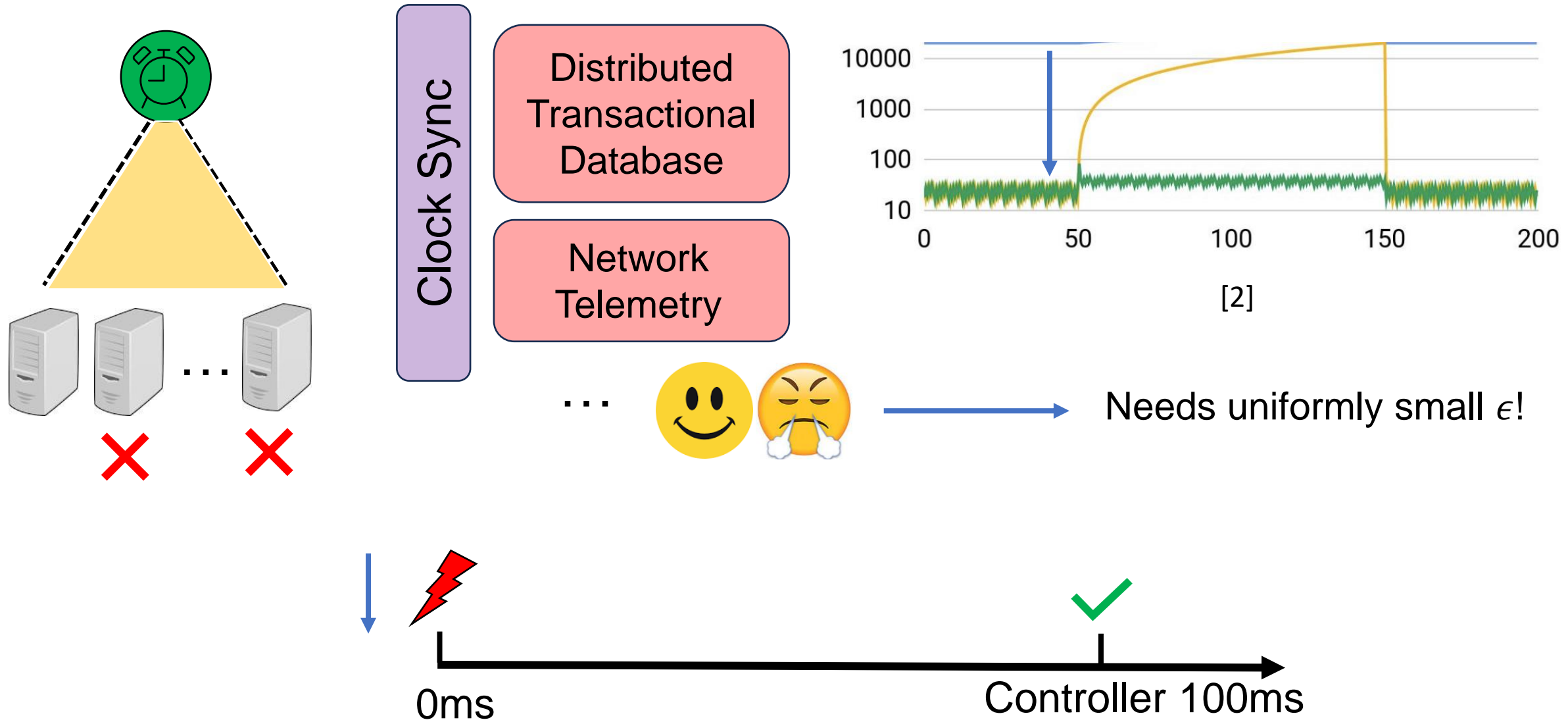
Global SDN Controller

With 1M multicast groups, Elmo [1] takes  $25ms + [\text{mcast group computation}]$  to handle network failures.

3ms completely no connection, and much more inefficient unicast (4.0x bandwidth consumption) for the rest of the time.



# Clock Synchronization



# Key idea: Data-planes

# and distributed algorithms

Distributed  
Pri-

**A Simple and Linear Time Randomized Algorithm for Computing Sparse Spanners in Weighted Graphs \***

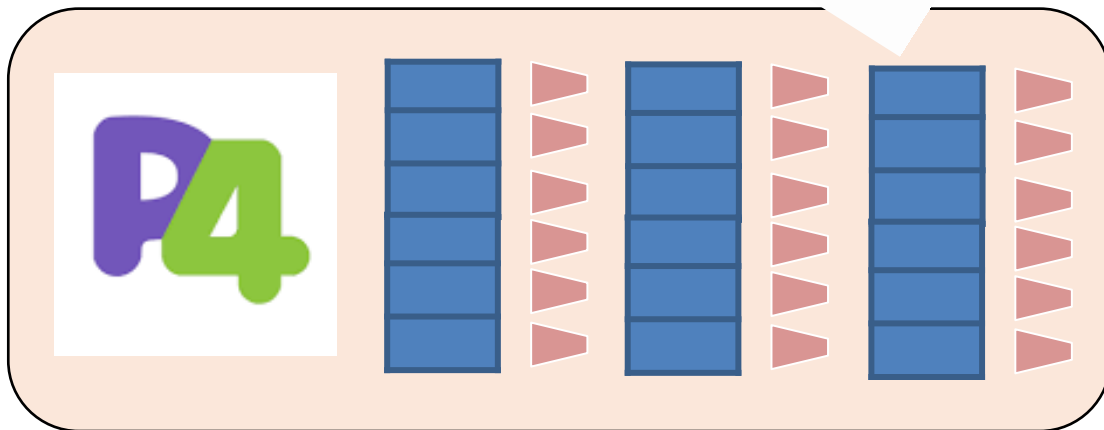
Surender Baswana<sup>†</sup>  
Max-Planck-Institut für Informatik,  
Stuhlsatzenhausweg 85,  
66123 Saarbrücken, Germany.  
Email: sbaswana@mpi-sb.mpg.de

Sandeep Sen<sup>‡</sup>  
Department of Comp. Sc. and Engg.,  
Indian Institute of Technology Delhi,  
Hauz Khas, New Delhi-110016, India.  
E-mail: ssen@cse.iitd.ernet.in

FACIS



network changes

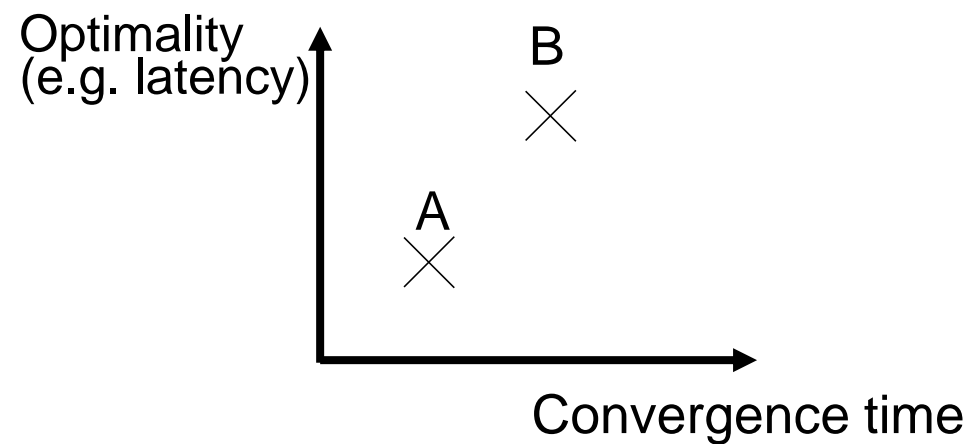


- **Benefits:** Fast reaction.
  - $O(100)$  ms  $\rightarrow$   $O(1)$  ms
- **Bridging** the gap between theory community and system in practice (bidirectional).

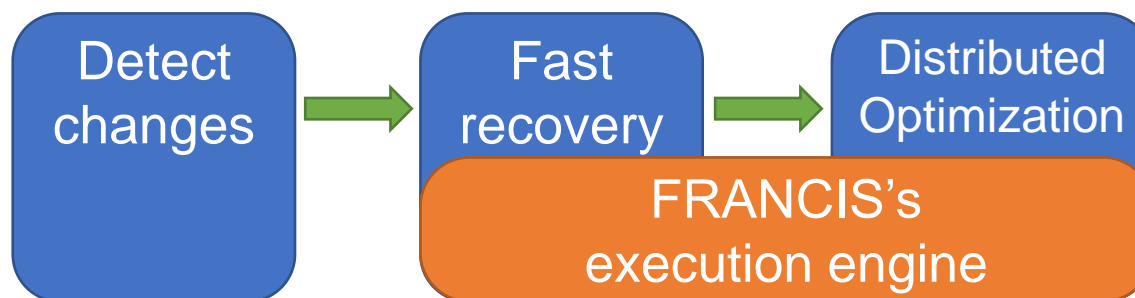


# The principle of choosing algorithms

- Fast recovery for reactivity, and dist optimization for optimality.

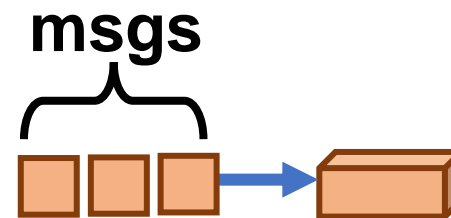


A: One spanning tree for <b>all</b> groups.	289 $\mu$ s	1.4 $\mu$ s
B: One mcast tree for <b>each</b> group.	1481 $\mu$ s	0.9 $\mu$ s
	Conv time	latency



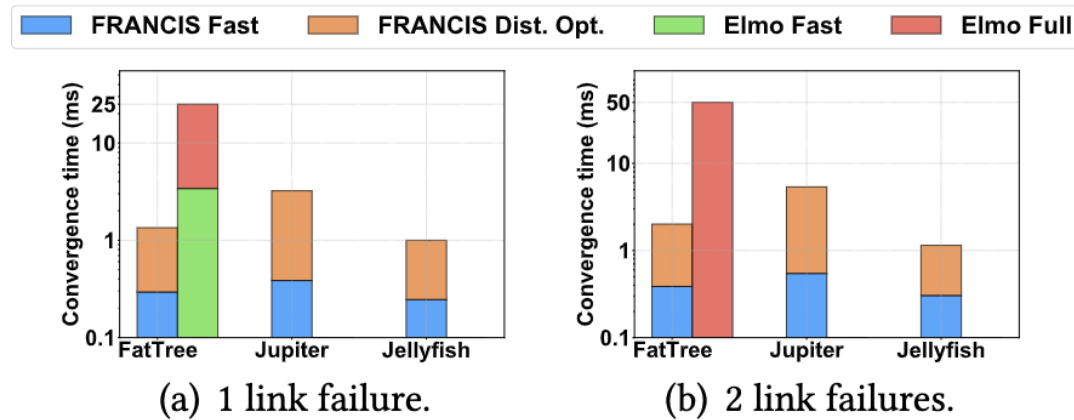
# FRANCIS's Design overview

Achieving generality	✓ Supporting both Async and Sync algorithms
	✓ $Bw \leq B, Mem \leq M$ $B$ and $M$ are user-specified.
Execution on switches	✓ Overcoming Tofino switch restrictions
	✓ Packet loss handling
	✓ Message-packing for reducing bandwidth consumption

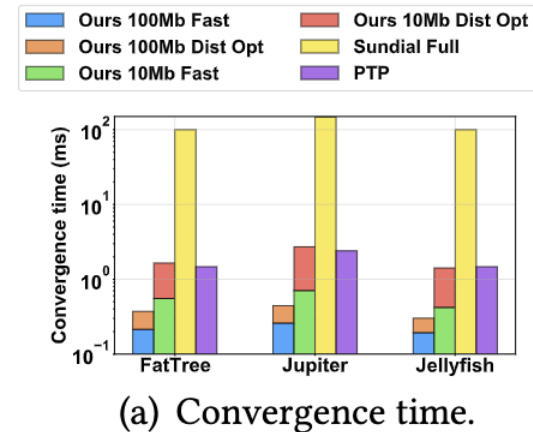


# Benefits: fast reaction

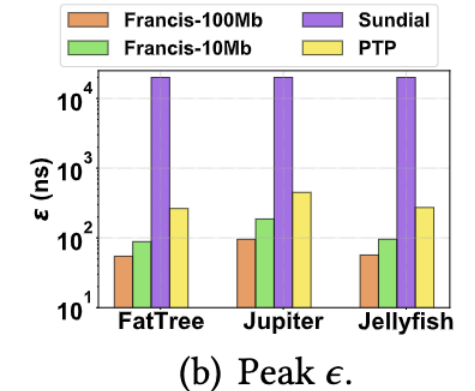
- FRANCIS's Data-plane-based approaches converge 1-2 orders of magnitude faster!



multicast



Clock synchronization



100Mbps bandwidth, 1.5% SRAM usage

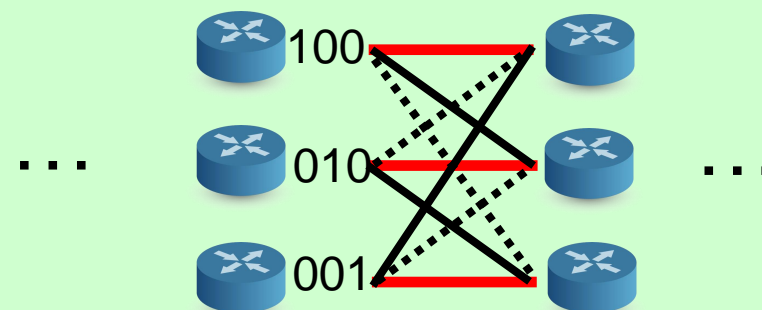
# Bridging the gap between theory and practice

## Theory

- Distributed spanner algorithm: find  $G' \subset G$  with minimum #edges such that  $d_{G'}(x, y) \leq \lambda d_G(x, y)$ .



## Source-routed multicast



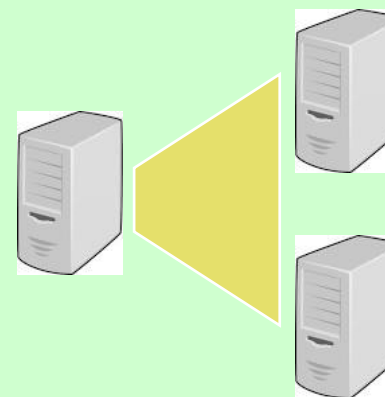
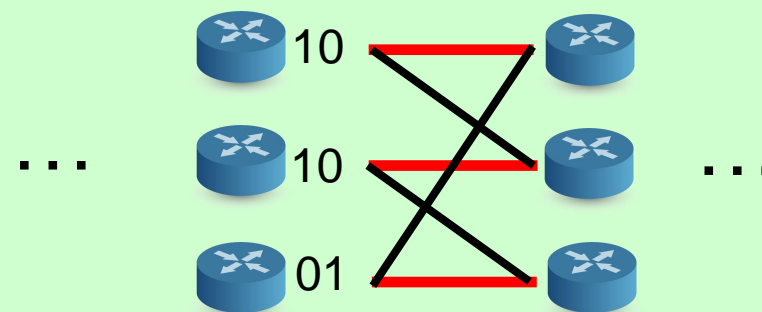
# Bridging the gap between theory and practice

## Theory

- Distributed spanner algorithm: find  $G' \subset G$  with minimum #edges such that  $d_{G'}(x, y) \leq \lambda d_G(x, y)$ .
- Minimum size shortest path tree.



## Source-routed multicast



- $\downarrow$  #Links, s.t.
- Lowest latency



Thank you for your attention!

Email: [wenchen.han.22@ucl.ac.uk](mailto:wenchen.han.22@ucl.ac.uk)

Q&A?