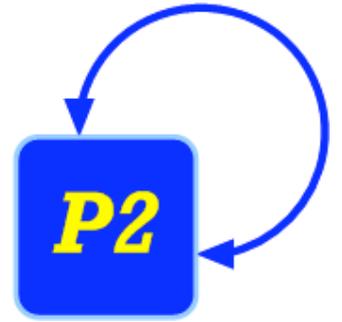# Finally, a Use for Componentized Transport Protocols

Tyson Condie, Joseph M. Hellerstein,
Petros Maniatis, Sean Rhea, Timothy Roscoe
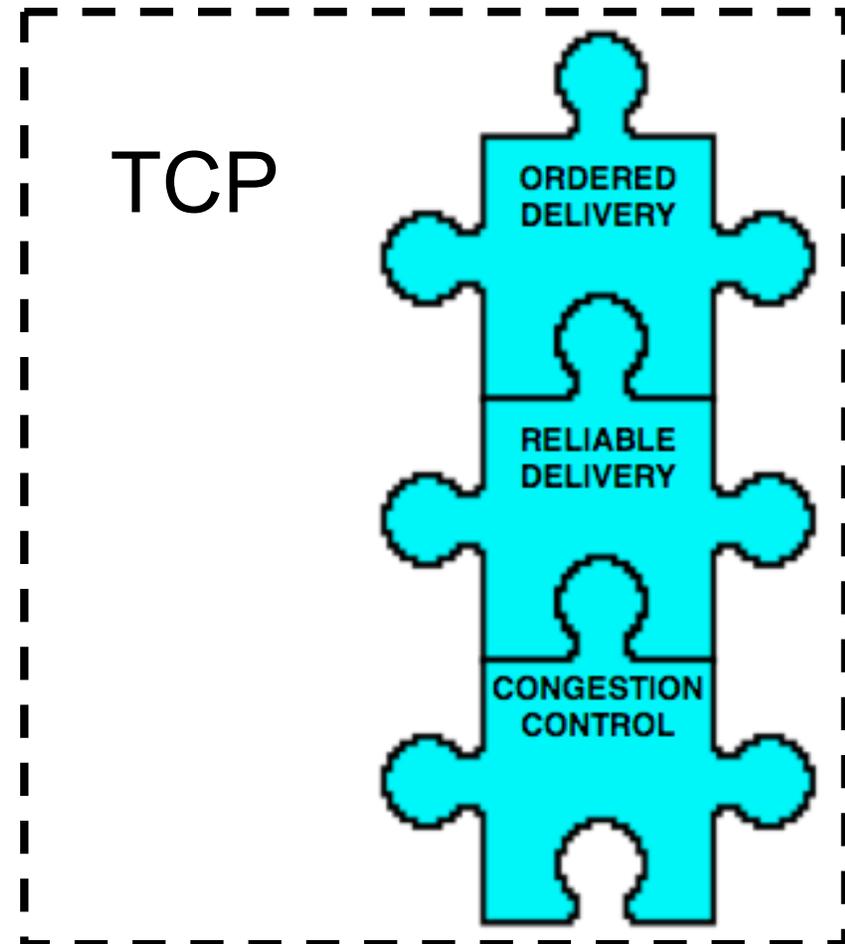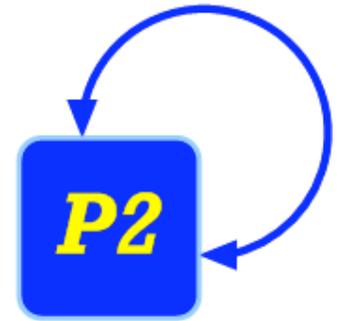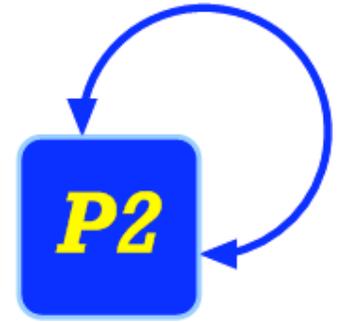U.C. Berkeley and Intel Research Berkeley

# Roadmap

- History
- What's different now?
- New uses for componentized transport protocols
- Conclusion
- Ongoing work

# Componentized Protocols?

**P2**

Decomposing transport protocols into a set of reusable building blocks that can be recomposed in different ways depending on application and network properties

TCP

ORDERED DELIVERY

RELIABLE DELIVERY

CONGESTION CONTROL

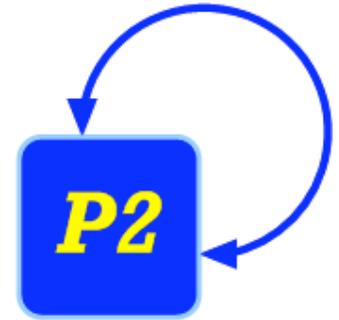# A Brief History of Componentized Protocols



- x-Kernel system
  - Export user-level protocol objects to compose into more general communication services

- Morpheus programming language
  - Object oriented programming support for protocol objects
  - Compiler time optimizations over generated protocols

- Prolac Protocol Language
  - Expression language for developing complete protocols

- Componentized protocols never caught on
  - Most applications were satisfied with point-to-point protocols
    - TCP, RTP, SCTP, DCCP, etc.

N. C. Hutchinson and L. L. Peterson. The x-kernel: An architecture for implementing network protocols. IEEE Trans. Software Eng., 17(1):64-76, 1991.

E. Kohler, M. F. Kaashoek, and D. R. Montgomery. A readable TCP in the Prolac protocol language. In Proc. ACM SIGCOMM, 1999.

4

# Why Now?

- Increasing popularity of overlay networks

  - DHT, BitTorrent, Akamai, Narada, Freenet

- Overlay networks have a broad design space

  - Nodes play the role of client, server, and router

- Most protocols today are tuned for point-to-point communications
  - Overlay requirements go beyond point-to-point model
  - Forces overlay programmers to develop their own handcrafted transport layer

- We have built a componentized protocol framework using a dataflow abstraction
  - Does it meet the set of opportunities and requirements of overlays?
  - Does it provide a programmer friendly framework?

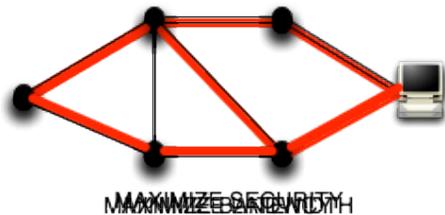*Every good work of software starts by scratching a developer's personal itch*

The Cathedral and the Bazaar, Raymond

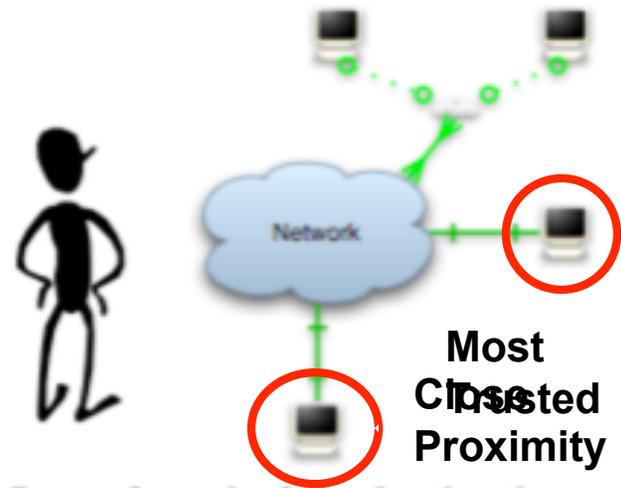# Application-level Routing Freedom

- A message need not be bound to a single destination or path
  - Several equivalent destinations
  - Several paths to get to a destination
- Fine grain application control over where a message is sent



Several paths to destination

MAXIMIZE SECURITY
MAXIMIZE BANDWIDTH

Most
Close
Trusted
Proximity

Several equivalent destinations

S. Rhea, D. Geels, T. Roscoe, and J. Kubiatowicz. Handling Churn in a DHT. In Proc. of the 2004 USENIX Annual Technical Conference, Boston, MA, June 2004.

A. Mislove and P. Druschel. Providing Administrative Control and Autonomy in Structured Peer-to-Peer Overlays. In Proc. of IPTPS 2004.

F. Dabek, B. Zhao, P. Druschel, J. Kubiatowicz, and I. Stoica. Towards a Common API for Structured Peer-to-Peer Overlays. In Proc. of IPTPS 2003.

S. Marti, P. Ganesan, and H. Garcia-Molina. SPROUT: P2P Routing with Social Networks. In Proc. of P2P&DB 2004.

P. Druschel, Y. Charlie Hu, J. Jannotti, M. Castro, A. Rowstron, and A.-M. Kermarrec. SplitStream: High-Bandwidth Multicast in Cooperative Environments. In Proc. of SOSP 2003.

F. Dabek, R. Cox, F. Kaashoek, and R. Morris. Vivaldi: A Decentralized Network Coordinate System. In Proc. of ACM SIGCOMM, Portland, OR, Aug. 2004.

S. Kamvar, M. Schlosser, and H. Garcia-Molina. The EigenTrust Algorithm for Reputation Management in P2P Networks. In Proc. WWW 2003.
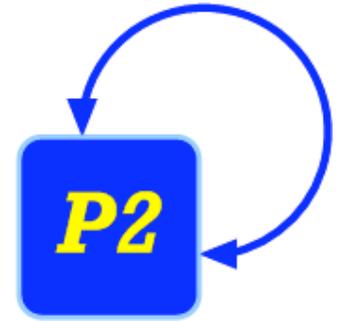
# Alternative Congestion Control

- **Per hop** congestion control
  - Recursive routing
    - UdpCC connection for each neighbor node

  S. Rhea, D. Geels, T. Roscoe, and J. Kubiatowicz. **Handling Churn in a DHT.** In *Proc. of the 2004 USENIX Technical Conference*, Boston, MA, USA, 2004.
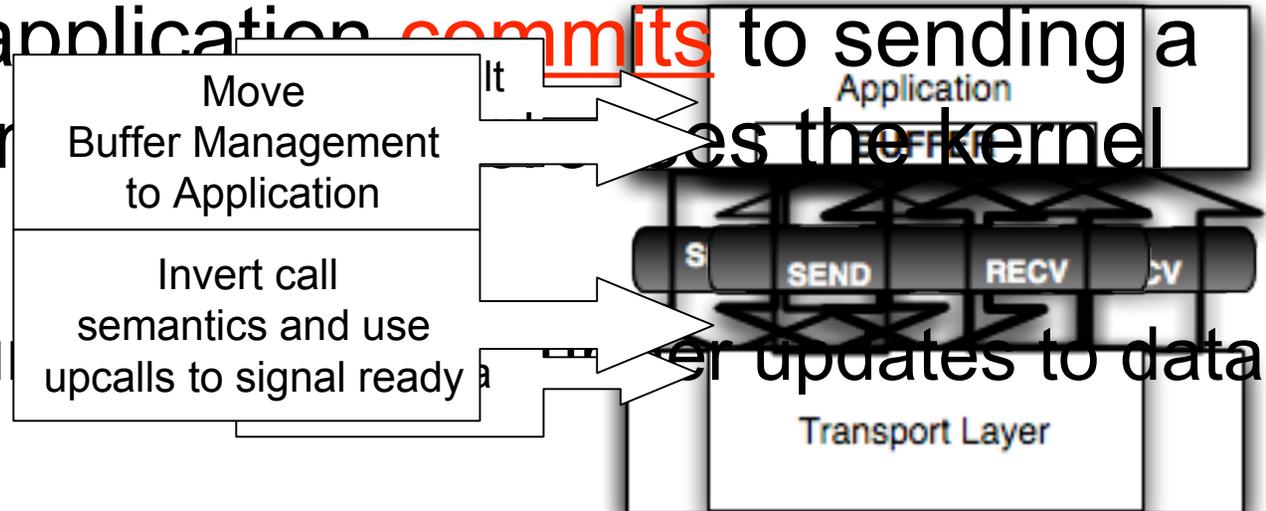
- **Aggregate** congestion control
  - Short lived connections
  - Connections that send very little traffic
  - Iterative routing (a la MIT Chord)
    - Next hop discovered during the lookup

  F. Dabek, J. Li, E. Sit, F. Kaashoek, R. Morris, and C. Blake. **Designing a DHT for low latency and high throughput**. In *Proc. NSDI*, 2004.
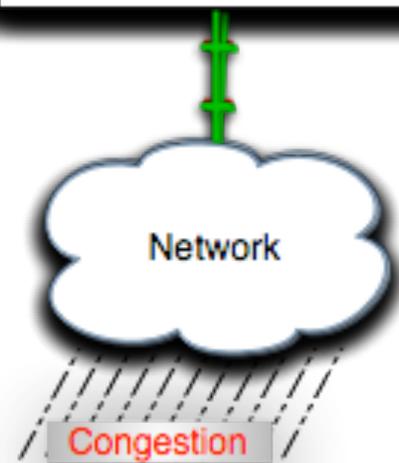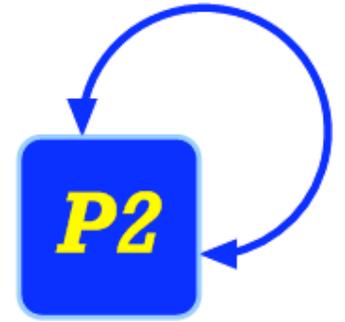
# Late Data Choice

- In TCP an application commits to sending a packet when it crosses the kernel boundary
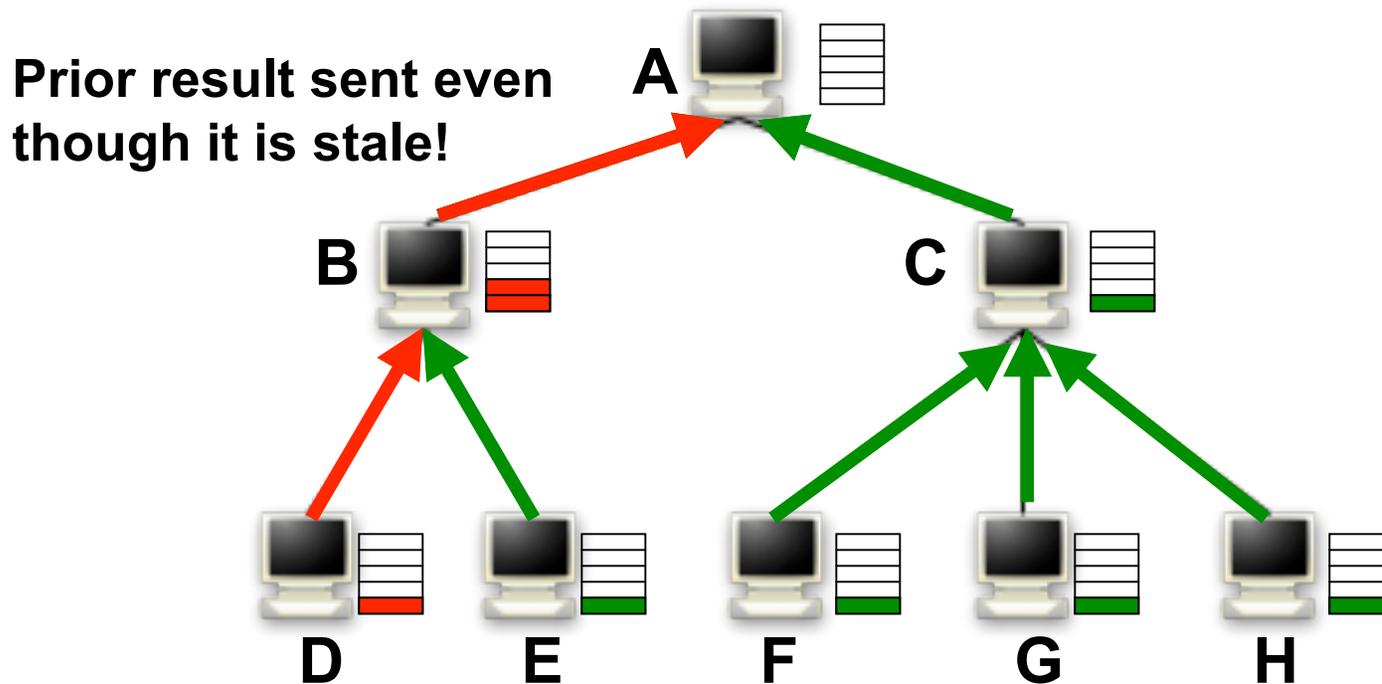  - Kernel boundary defines latest point for updates to data



**Move Buffer Management to Application**

**Invert call semantics and use upcalls to signal ready**

Application

SEND    RECV

Transport Layer

Network

Congestion

P2

David D. Clark, The Structuring of systems using upcalls, ACM SIGOPS Operating Systems Review, v.19 n.5, p.171-180, Dec. 1-4, 1985

Moogul Brakmo Peterson DF Towards Operating System Predictability, SIGCOMM Comput. Commun. Rev. 34(1) pp.106 2004
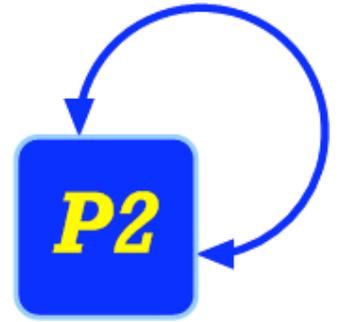
8

# Late Data Choice

- Application sensitive to stale results
- Transmission costs are high
  - Late data choice ensures the most up-to-date computation is sent



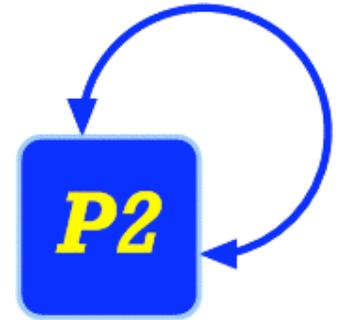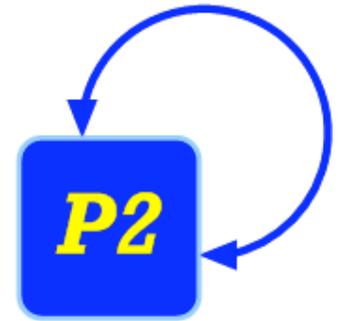**Prior result sent even though it is stale!**

# Roadmap

- History
- What's different now?
- New uses for componentized transport protocols
- Conclusion
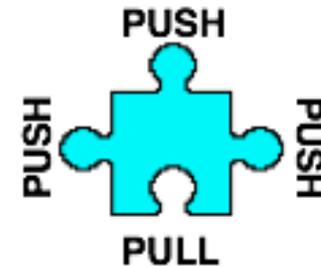- Ongoing work

# Benefits to Componentized Protocols

- Fine grain protocol modifications
  - Alternative congestion control
  - Message/packet level reliable delivery
  - Custom packet scheduling algorithms
- Transport layer more knowledge visible
  - Late data choice
  - Transport state can aid failure detection, replica selection, load balancing decisions, etc.
- Encode domain knowledge in the transport layer
  - Overlay routing logic
  - Message semantics

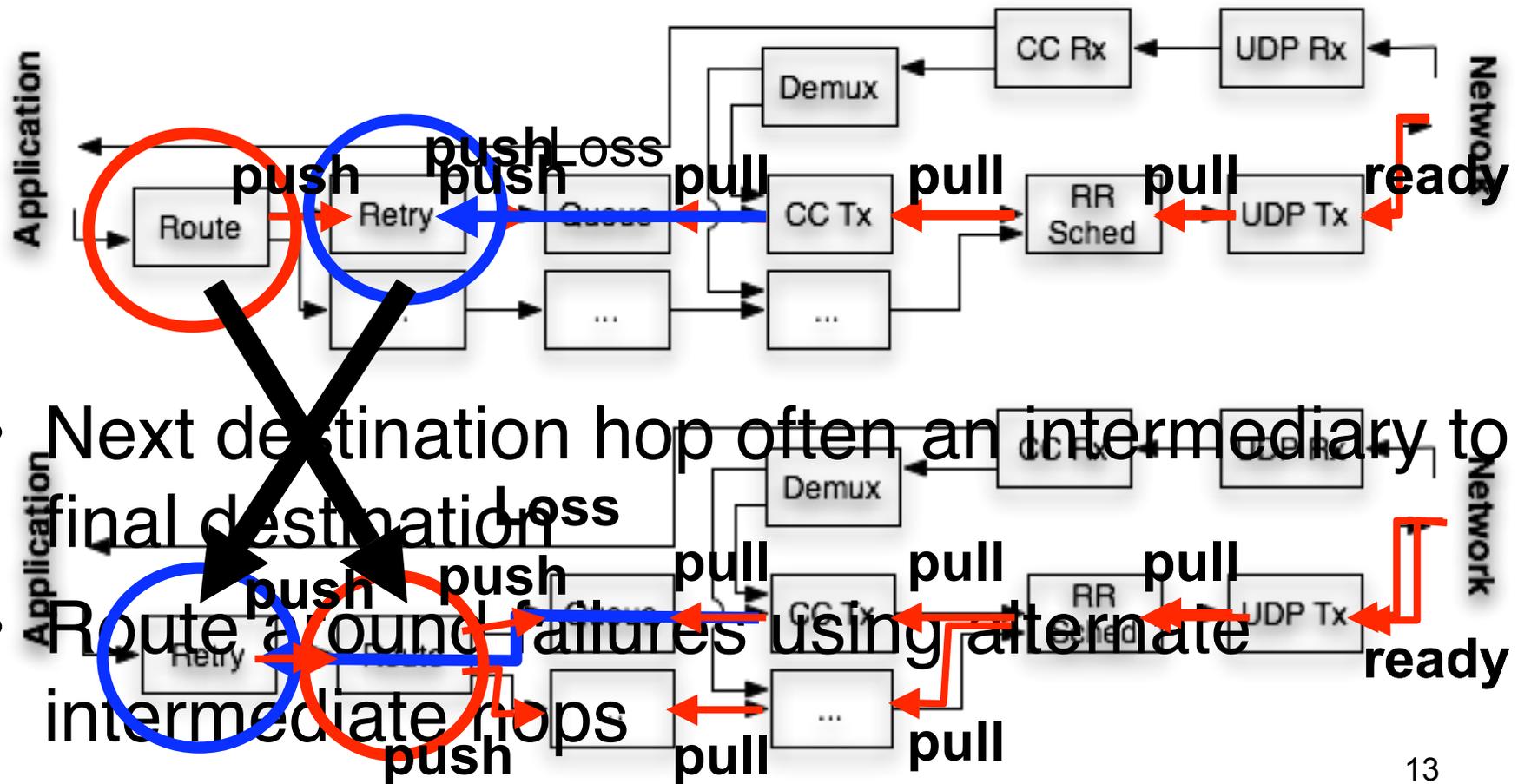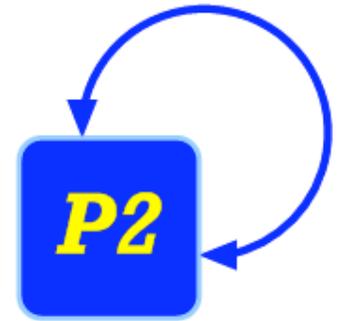# Componentized Protocols using a Dataflow Abstraction

- ## Graph model places elements at the vertices
  - Elements abstract code into modular units that perform a specific task
  - Elements export push or pull **interface**

- ## Graph structure orders data transformations
  - Traditional protocols follow stack ordering
  - Dataflow more general
    - Protocol semantics encoded into the graph structure

E. Kohler, R. Morris, B. Chen, J. Jannotti, and M. F. Kaashoek. **The Click modular router.** ACM Trans. Comput. Syst., 18(3):263-297, 2000.
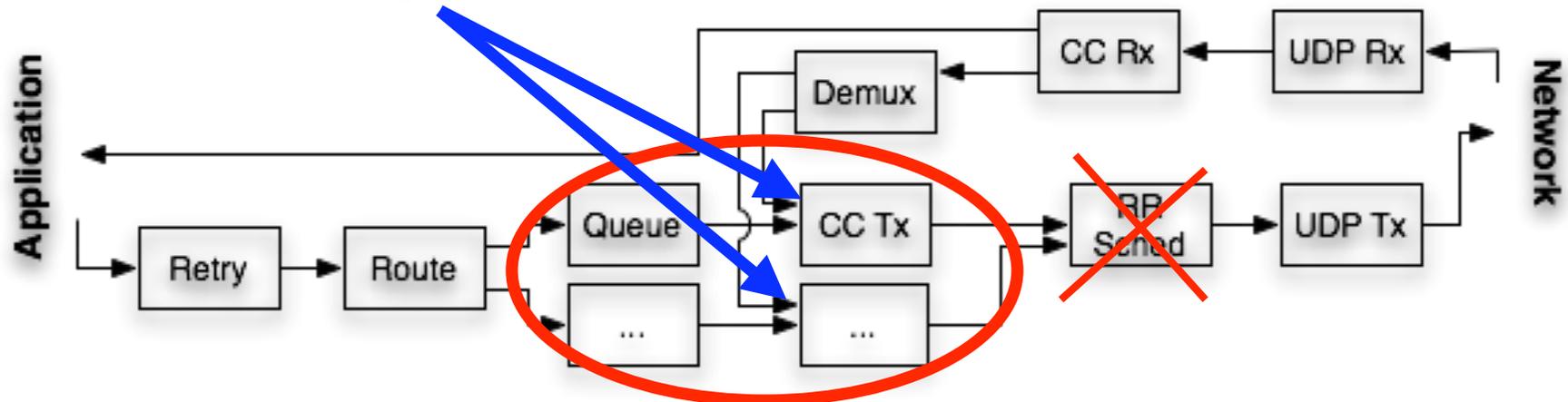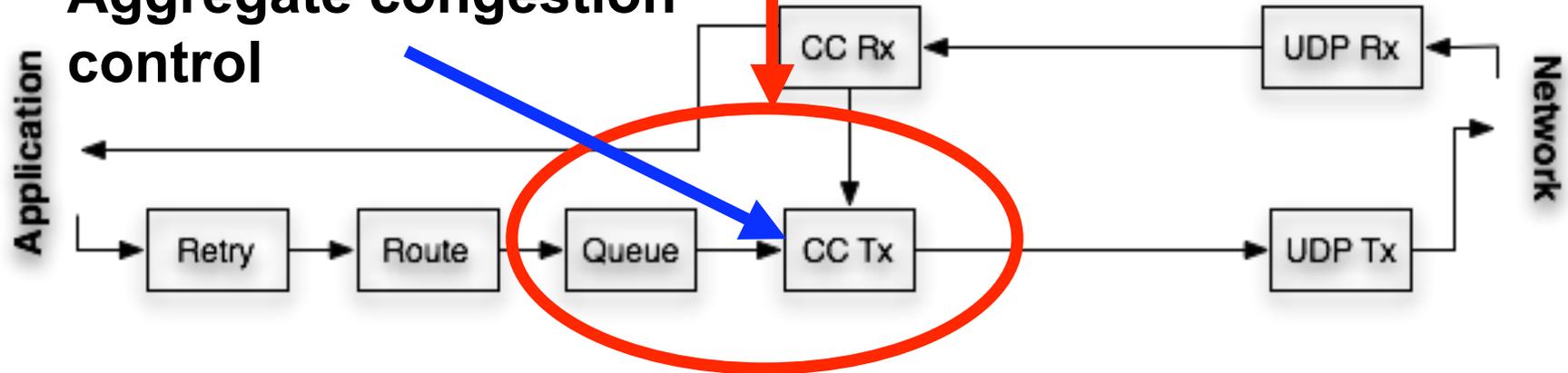
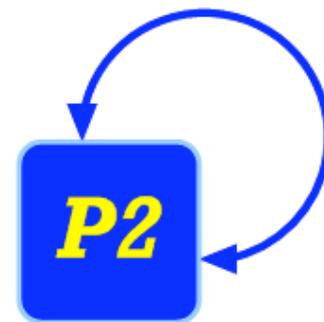# Application-level Routing Freedom

**P2**



- Next destination hop often an intermediary to final destination

- Route around failures using alternate intermediate hops

13

# Aggregate Congestion Control

**Per hop congestion control**

**Aggregate congestion control**

# Late Data Choice

- Network is asynchronous

**Buffer beyond application control**
  - Message arrival and transmit times are unpredictable
- Push based dataflow accepts whatever data it is given
  - Follows semantics of network receive
- Pull based dataflow awaits a signal before releasing its data
  - Follows semantics of network send

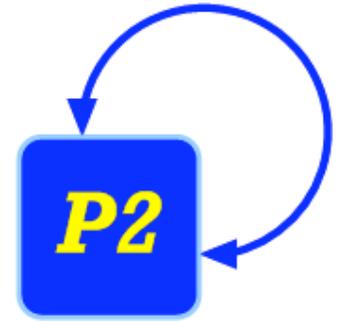**Better queuing properties!**

**Signaling transparent**

Network receive entirely **PUSH** based

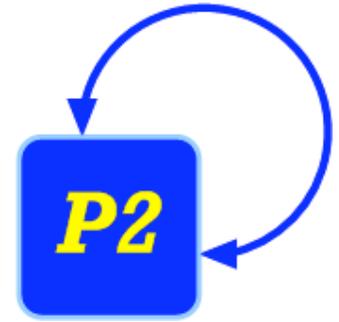Network send entirely **PULL** based

15

# We can do all this in P2

- P2: A query processor for constructing overlays
  - Uses a declarative language for specifying queries that describe overlay properties/invariants
  - Queries compiled into a dataflow graph
- P2 dataflow model extends into network stack
  - Satisfies our transport layer needs for building overlays
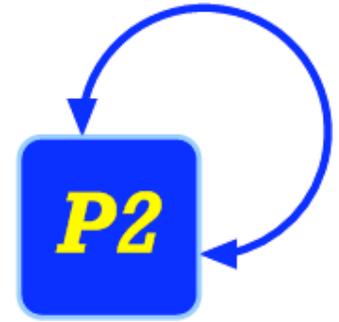  - Blurs boundary between application and transport

B. T. Loo, T. Condie, J. M. Hellerstein, P. Maniatis, T. Roscoe, and I. Stoica. **Implementing declarative overlays.** In *Proc. ACM SOSP 2005*.
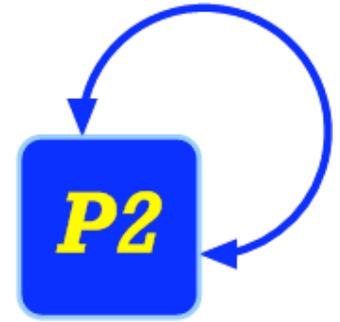
# Conclusion

- Overlays offer many new design decisions
  - Functionality requirements go beyond the scope of current monolithic transport services
  - Requirements well suited to componentized protocols

- Componentized protocols erase the functionality boundary
  - Can encode the application and transport layer with the right set of features

- Dataflow is an instance of componentized protocols
  - Flexible glue layer between network and application
  - Code reuse through graph modification

# *Ongoing Work*

- Declarative language for transport layer
  - Translate high level invariants into supporting dataflow(s)
- Automatic static generation of dataflow graphs
  - Each semantically equivalent dataflow can offer certain application and network tradeoffs
  - Cost model chooses an optimal dataflow to install
- Runtime reconfiguration / reoptimization
  - What kinds of modifications and how are they triggered?
  - What kinds of statistics would aid in this effort?

# Thank You!

For more information please visit our website
http://p2.cs.berkeley.edu/