# Paper Review of FlowBender

Liang Wu*
liangwu@uwm.edu
University of Wisconsin, Milwaukee
Milwaukee, Wisconsin, USA

## 1 Introduction

This paper explores the topic of datacenter networks, focusing on the challenges in managing datacenter traffic. Datacenter networks provide multiple paths for traffic between machines, and effective load balancing across these paths is crucial for applications sensitive to latency and throughput. When traffic is not evenly distributed, some paths can become congested while others are underutilized, negatively impacting network throughput and increasing tail latency for flows. High tail latency in datacenter applications directly affects the end-user experience because user responses are aggregated from the results of thousands of servers, making the tail latency of individual flows a bottleneck for response time.

A commonly used load-balancing mechanism in datacenter networks is Equal Cost Multiple Path (ECMP) forwarding. This method randomly assigns flows to different paths by hashing specific fields in the packet headers, ensuring that packets within a flow follow the same path. Although ECMP is widely adopted because it works seamlessly with off-the-shelf switches and standard TCP/IP stacks, it has limitations in handling datacenter traffic effectively. In real-world applications, datacenter traffic often includes a small number of long flows that contribute a significant portion of the total traffic load. These long flows lack sufficient entropy for uniform load distribution, causing multiple long flows to collide on the same paths, leading to prolonged congestion, while other paths remain underutilized.

Several recent approaches have attempted to address ECMP's shortcomings. MPTCP, an end-to-end scheme, dynamically adapts to network conditions, unlike the static nature of ECMP. However, MPTCP introduces significant overhead in managing packet reordering. Random Packet Scatter (RPS) is a static solution that also experiences high levels of packet reordering. Other methods, like DeTail, require hardware-level changes, which are costly and hinder continuous deployment.

In this paper, the researchers propose FlowBender, a dynamic alternative to ECMP that minimizes packet reordering. FlowBender addresses ECMP's core limitation: the flow-to-path mapping is static; there is no consideration for re-routing if the current path is congested or even broken. FlowBender's design is based on the observation that a small number of long flows contribute a large fraction of network load. To mitigate congestion, FlowBender selectively reroutes flows only during congestion or route failures, thus avoiding unnecessary packet reordering costs.

FlowBender offers several advantages: it requires no hardware changes to switch silicon, involves only 50 lines of kernel code modifications, and can be implemented with simple reconfiguration of ECMP hash functions (through a few commands). It significantly outperforms ECMP and matches the performance of more complex schemes. Additionally, FlowBender incorporates robust end-to-end congestion notifications (ECN) and failure signals (timeouts). It reroutes at the granularity of Round-Trip Time (RTT) and quickly recovers from link failures within a timeout period.

The researchers evaluated FlowBender through both simulations and real-world implementation. Their results indicate that FlowBender reduces tail latency by over 40% on average for large flows compared to ECMP. In large-scale ns-3 simulations using datacenter-like workloads, FlowBender decreases the mean and tail latency of all-to-all workloads by 73% and 93%, respectively, compared to ECMP, while performing within 2% of costlier alternatives like DeTail and RPS. For storage-type workloads that produce incast traffic, FlowBender enables jobs to complete in half to a quarter of the time compared to ECMP, closely matching the performance of more complex schemes like DeTail and RPS (within 2% in average completion time and tail latency). Overall, FlowBender demonstrates the effectiveness of a non-intrusive, flow-level load balancing solution based on the classic end-to-end principle.

## 2 Approach Overview

As noted in the introduction, this paper examines the challenges of load balancing in data center networks,

emphasizing the importance of evenly distributing traffic across multiple paths to prevent congestion. This section will delve into the primary issues faced by data centers in detail and discuss current solutions, along with their strengths and weaknesses.

## 2.1 Main Problem

ECMP is the standard mechanism used in datacenters for load balancing data across multiple paths. In ECMP, the switch calculates a hash function based on the packet's header fields. The resulting hash value determines the port used to forward the packet. ECMP assigns flows to paths based on this static hash value derived from each flow's identifying fields. Once a flow is assigned to a path, it consistently uses that same path, with no dynamic adjustment based on traffic conditions. For short-lived flows, this design works effectively. However, issues arise with long-lived flows (flows that persist for an extended period), as they can lead to collisions. When two flows generate the same hash value, they collide and are assigned to the same path, regardless of other available paths. If these flows are long-lived, they will continuously occupy and congest that path. Because ECMP's flow-to-path assignment is static, it cannot redistribute or rebalance these long flows across other, less congested paths.

## 2.2 Existing Solutions

Many solutions aim to achieve better load balancing in datacenters and address the issues associated with ECMP. MPTCP (Multipath TCP) splits a single TCP flow into multiple subflows, each with a different port number. Since port numbers are part of the ECMP hashing function, these distinct subflows are hashed to different paths in the network, distributing the traffic load more evenly. By splitting traffic across multiple paths, MPTCP helps avoid the problem that arises when ECMP assigns long flows to the same path, allowing for dynamic traffic distribution based on network conditions. In MPTCP, the sender has additional responsibilities: it monitors metrics such as Round-Trip Time (RTT), throughput, and packet loss for each subflow and prioritizes paths with better performance. The receiver, on the other hand, requires extra logic to reassemble data from multiple subflows back into the original stream, which adds complexity. MPTCP has demonstrated improvements over ECMP for larger flows (those over 70 KB) as it effectively balances load across multiple paths, avoiding the limitations of ECMP's static path assignment for long flows. However, MPTCP has some

drawbacks. The congestion control algorithm on the sender side and the data reassembly logic on the receiver side are complex, adding over 10,000 lines of kernel code. This complexity can slow down performance and increase the difficulty of implementation and maintenance. Additionally, MPTCP incurs significant CPU overhead, consuming around 40% of CPU on clients and 10% on servers—considered excessive for datacenter environments.

On the other hand, Random Packet Spraying (RPS) and DeTail shift the responsibility for load balancing to the network switches. RPS operates at the packet level rather than the flow level, meaning each packet within a flow can take a different path through the network. For every packet, the switch randomly selects an output port from a set of available ports, independently forwarding each packet along different paths. RPS performs well in symmetric topologies like fat-tree structures, where paths between devices are balanced and equal. In such cases, packets reach their destinations efficiently, even when they are randomly sprayed across paths. However, when the network experiences asymmetries—due to link failures or incremental deployments (where new equipment is added over time)—RPS can lead to significant throughput inefficiencies, as it may overload some paths while underutilizing others.

DeTail overcomes RPS's limitations by combining Priority Flow Control (PFC) and adaptive routing to enhance tail latency and network utilization. PFC ensures lossless packet forwarding by preventing packet drops. However, PFC comes with challenges, such as the risk of network deadlock (when packets are indefinitely stuck in the network) and head-of-line (HOL) blocking within the same priority class, which can delay traffic. DeTail uses a form of adaptive routing where each packet chooses the least congested eligible port at each switch. This strategy helps reduce congestion by dynamically balancing traffic based on real-time network conditions. However, adaptive routing also presents challenges. It requires custom switch silicon to monitor queue lengths in real-time and to forward packets accordingly. This specialized hardware is expensive and has long production cycles, making rapid deployment challenging. Another issue is that, since DeTail performs load balancing at the per-packet level, packets can take different paths, leading to out-of-order delivery. This disrupts the TCP flow and disables TCP's fast retransmit feature, which is used to quickly resend lost packets. Lastly, because DeTail operates at the link

layer, it cannot manage failures on paths that have only one way to reach the destination. If a link fails in a part of the network where there is no alternate path (a deterministic part of the path), DeTail cannot reroute around the failure. It lacks mechanisms to switch the flow to an entirely different route that bypasses the failed link.

## 2.3 FlowBender Design

The core design principle of FlowBender adopts a simple approach that aligns with the requirements of modern data centers. FlowBender avoids introducing additional hardware overhead and minimizes the risk of out-of-order packet delivery. Its strategy focuses on transmitting a flow along a single path and rerouting the flow only when congestion occurs. The two primary components of the FlowBender design are congestion detection and flow rerouting.

To detect congestion, FlowBender employs DCTCP (Data Center Transmission Control Protocol), which is part of the ECN (Explicit Congestion Notification) congestion notification scheme—a standard supported by contemporary switches. In DCTCP-style marking, a switch marks packets as congested when the queue size exceeds a predefined threshold. The sender then tracks the number of marked acknowledgments (ACKs) received during each round-trip time (RTT). If the number of marked ACKs exceeds a set threshold for a flow, the system interprets this as an indication of congestion, triggering rerouting.

For rerouting, FlowBender leverages the existing ECMP hash engine mechanism. ECMP in switches computes a hash value based on specific fields in a packet header, which uniquely identifies a network connection. FlowBender utilize the flexibility of hash engines, which can be configured to hash on fields beyond the typical connection identifiers. Specifically, FlowBender uses a flexible field in the packet header to compute the hash value. Each TCP socket maintains a value, V, that is inserted into the hashing function for every outgoing packet. When congestion is detected, V is updated, causing packets to be rerouted. This approach allows the network to adapt to congestion dynamically without requiring coordination or pre-agreement between the sender and the receiver. FlowBender selects flexible fields, such as VLAN ID or TTL, that can serve as hash inputs without disrupting packet delivery or requiring additional negotiation. These fields are ideal for rerouting because they do not affect the functionality of the data center network. In conclusion, FlowBender enables rerouting by having each TCP flow independently

maintain a value V and track the per-RTT fraction of marked acknowledgments (F). When F surpasses a predefined threshold (T), the flow detects congestion on the current path. To address this, the system modifies the value V in the flexible field, effectively rerouting the flow to a less congested path. This design ensures compatibility with existing data center infrastructure while enhancing network responsiveness to congestion.

# 3 Strength and Weakness

## 3.1 Strength

FlowBender demonstrates strong performance compared to DeTail and RPS. The researcher evaluated the performance of ECMP, RPS, DeTail, and FlowBender using traffic patterns representative of large-scale online services, such as web search. Across all traffic loads, DeTail, FlowBender, and RPS significantly outperformed ECMP. While RPS and FlowBender generally achieved similar levels of performance, FlowBender stands out because it does not require any hardware modifications to switches and is not strictly dependent on symmetric topologies, unlike DeTail and RPS.

FlowBender is straightforward to implement, with a design that requires minimal modifications. Its complete implementation involves approximately 50 lines of kernel code on the hosts and only 5 lines of configuration code on the switches, making it both lightweight and practical for deployment.

FlowBender is highly effective in minimizing packet reordering by employing a careful rerouting strategy. It waits for at least one Round-Trip Time (RTT) before rerouting a flow, which helps avoid excessive packet reordering. Only packets transmitted immediately before and after the reroute may arrive out of order, resulting in significantly less reordering compared to methods like DeTail or RPS (Random Packet Spraying). Additionally, FlowBender effectively handles false alarms, such as short-term traffic bursts (bursty marking), which might otherwise trigger unnecessary rerouting. By waiting for one RTT, FlowBender avoids reacting to these transient signals, ensuring rerouting decisions are stable and minimally disruptive.This approach allows FlowBender to reroute with minimal impact on packet order, maintaining robust performance and efficiency under challenging network conditions, unlike systems such as DeTail and RPS, which experience frequent and costly reordering. The researcher analyzed

the occurrence of out-of-order packets in all simulations. With FlowBender, the probability of a packet arriving out of order was only 0.006% higher than ECMP's. In comparison, DeTail experienced more than 97.9% of the out-of-order packet deliveries observed with RPS, further highlighting FlowBender's efficiency and robustness.

### 3.2 Weakness

FlowBender operates on top of ECMP, which provides the initial static flow-to-path mapping based on a hashing mechanism. ECMP itself is a static load-balancing approach. FlowBender introduces dynamism by modifying the flexible hashing field, effectively recalculating the ECMP hash to reassign the flow to a different path. ECMP assumes paths have equal cost and symmetric characteristics, which is often not true in real-world networks. This static assumption undermines the performance of FlowBender when network topologies are highly asymmetric.

In extreme congestion scenarios, FlowBender may continuously switch between paths if newly selected paths remain congested. This behavior could lead to unnecessary rerouting and packet reordering, impacting overall performance.

## 4 Inspiration

An enhancement to FlowBender could involve integrating machine learning models at switches to predict congestion and optimize flow assignments proactively. These models would analyze historical traffic data to identify potential congestion links, enabling high-demand flows to be directed to less congested paths. This approach would reduce the reliance on reactive rerouting, complementing FlowBender's existing mechanism and minimizing disruptions. Additionally, the models could address ECMP's limitations in handling asymmetric topologies by applying weighted probabilities that account for varying link capacities. To ensure adaptability, the models could be continuously retrained with real-time traffic data and redeployed, allowing the system to dynamically adjust to changing network conditions

By combining proactive machine learning optimization with FlowBender's congestion-driven rerouting, the system would achieve a smarter, more efficient load-balancing approach, improving performance in dynamic data center environments.

## 5 Conclusion

FlowBender offers an effective solution for load balancing in data center networks by dynamically rerouting flows in response to congestion. Its lightweight design requires minimal modifications, making it practical for deployment without hardware upgrades. FlowBender stands out for its strong performance compared to more expensive and complex solutions, as well as its ability to minimize packet reordering. It significantly reduces mean and tail latency compared to ECMP while achieving performance comparable to costly or less practical schemes like DeTail and RPS. Despite these advantages, FlowBender's reliance on ECMP introduces limitations, such as inefficiencies in asymmetric topologies and susceptibility to frequent rerouting during extreme congestion. Overall, FlowBender keeps a balance between simplicity and performance, significantly reducing tail latency and improving flow transmission time. It serves as a practical alternative method to complex schemes like MPTCP and DeTail, demonstrating strong potential for future enhancements.

## References

[1] Abdul Kabbani, Balajee Vamanan, Jahangir Hasan, and Fabien Duchene. 2014. FlowBender: Flow-level Adaptive Routing for Improved Latency and Throughput in Datacenter Networks. In Proceedings of the 10th ACM International on Conference on emerging Networking Experiments and Technologies (CoNEXT '14). Association for Computing Machinery, New York, NY, USA, 149–160. https://doi.org/10.1145/2674005.2674985