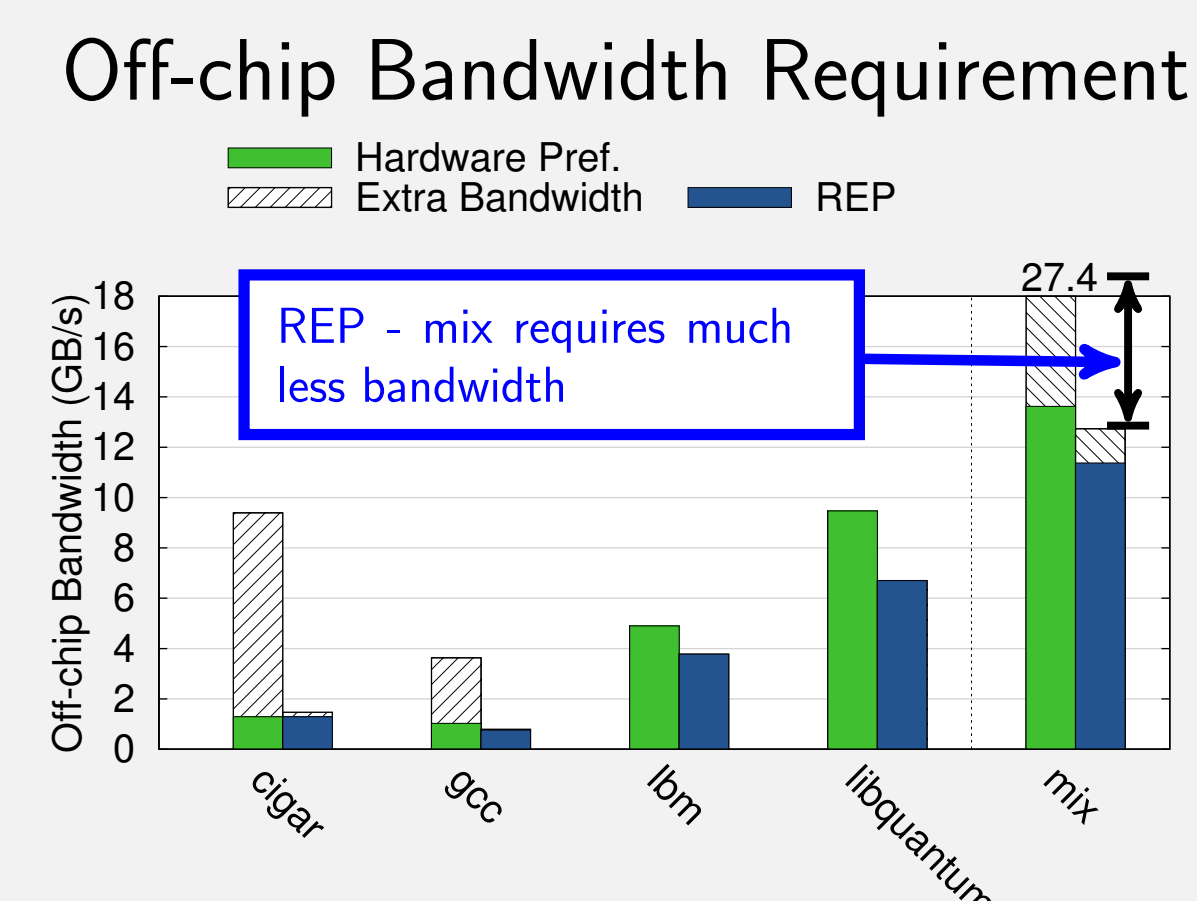
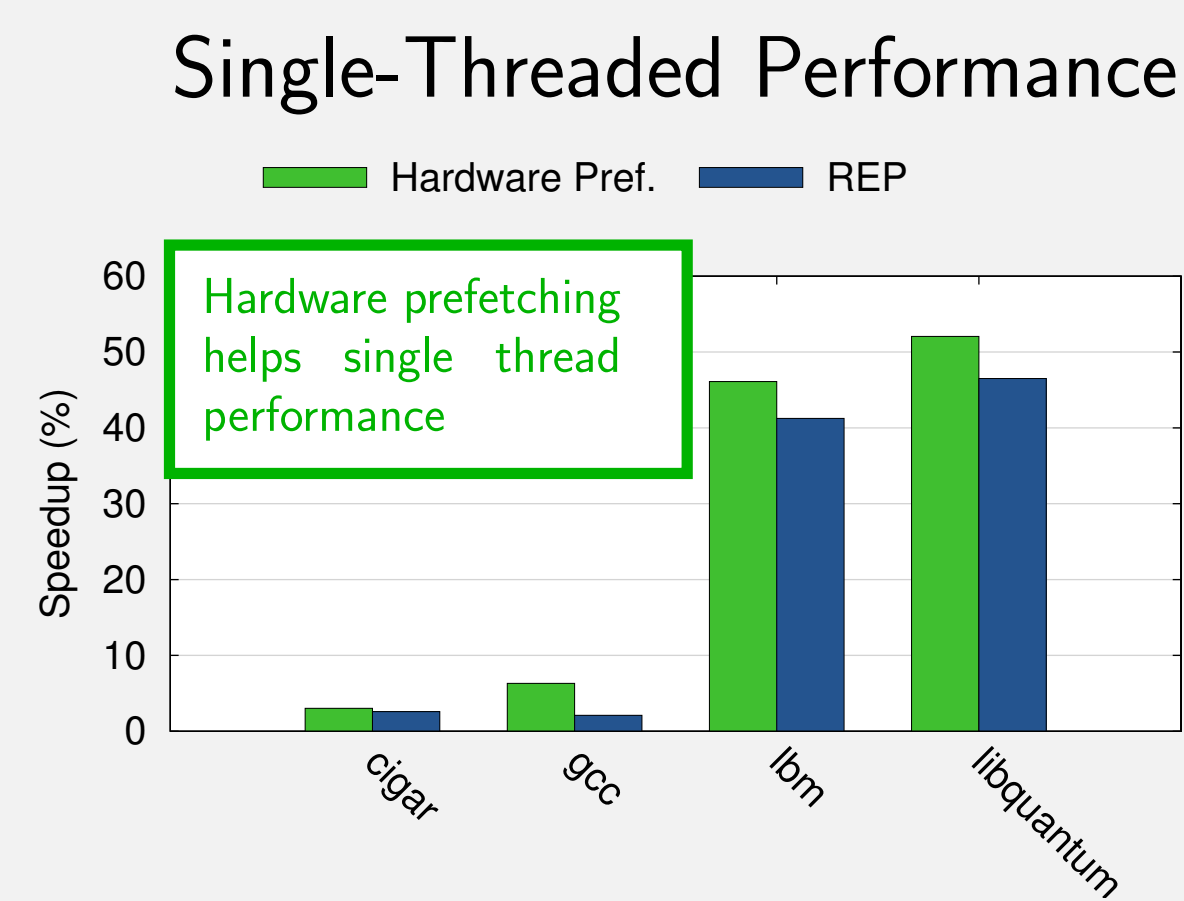
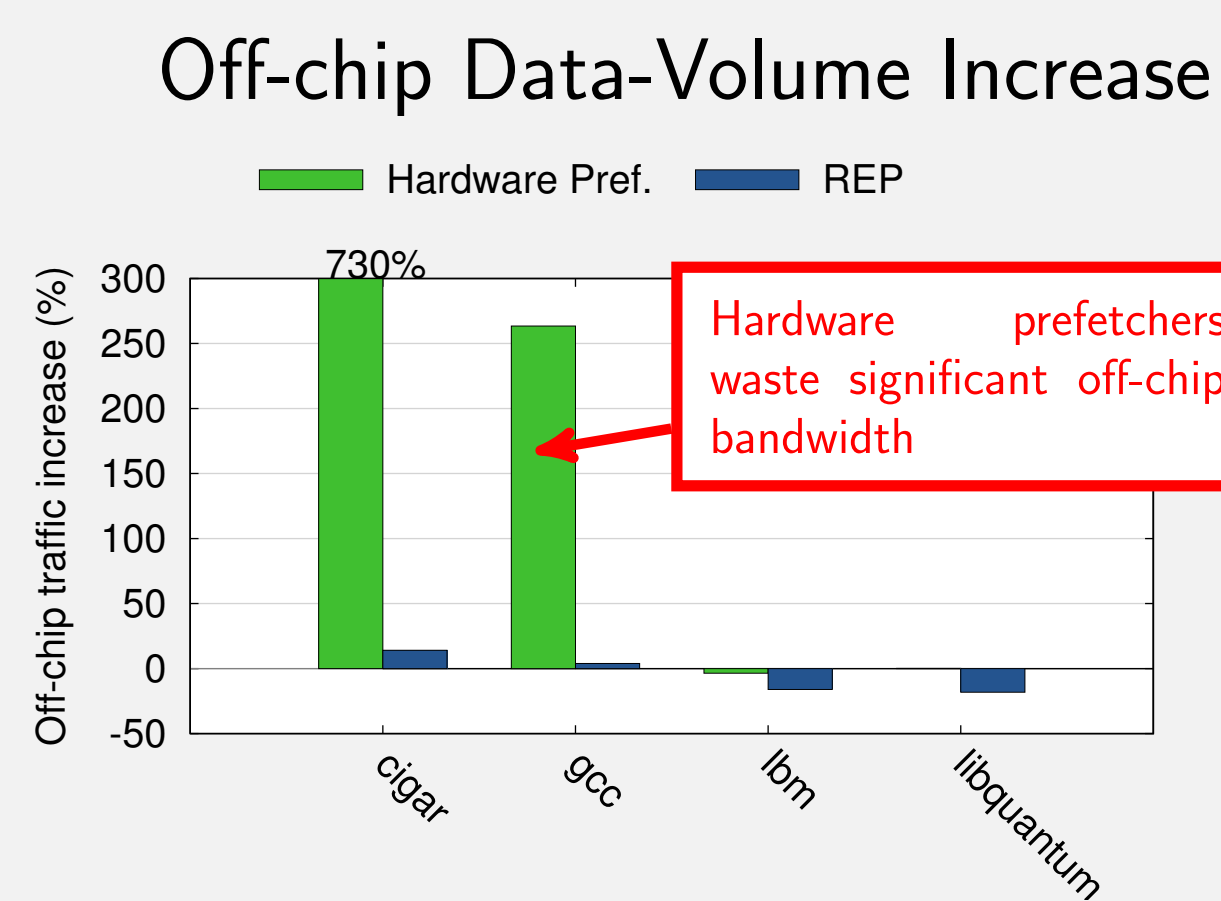
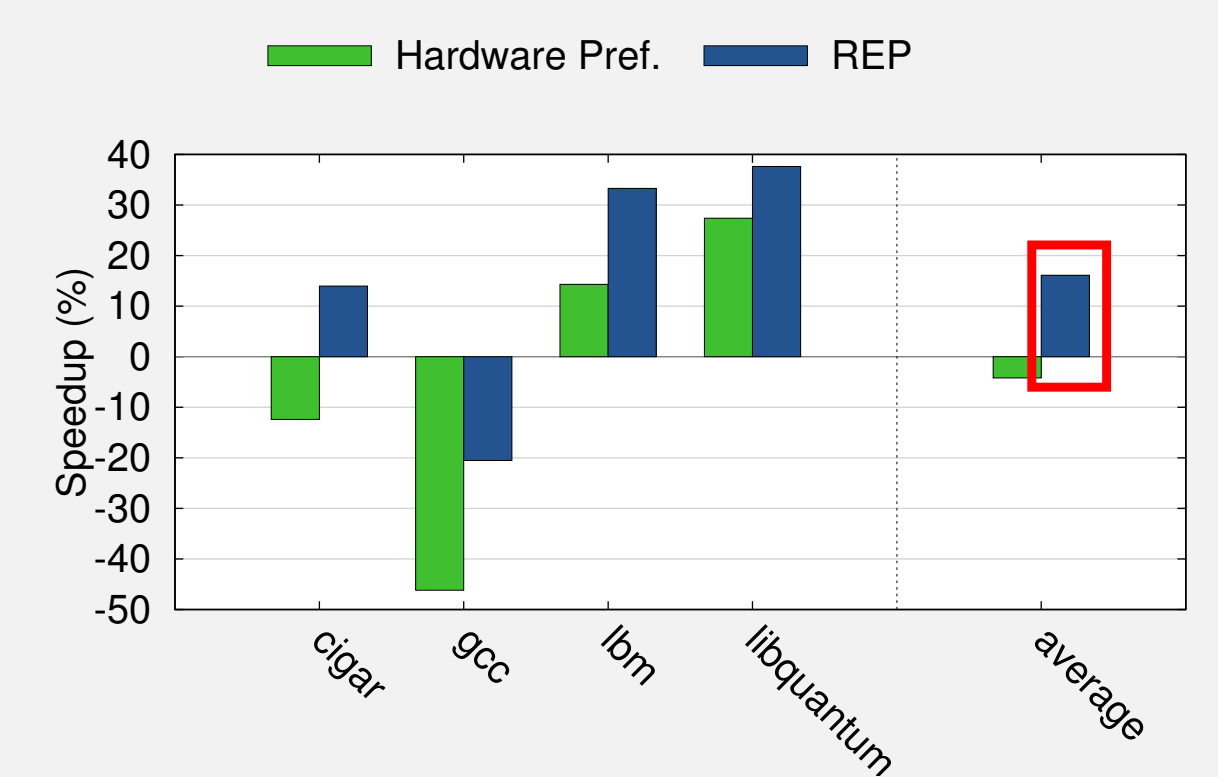


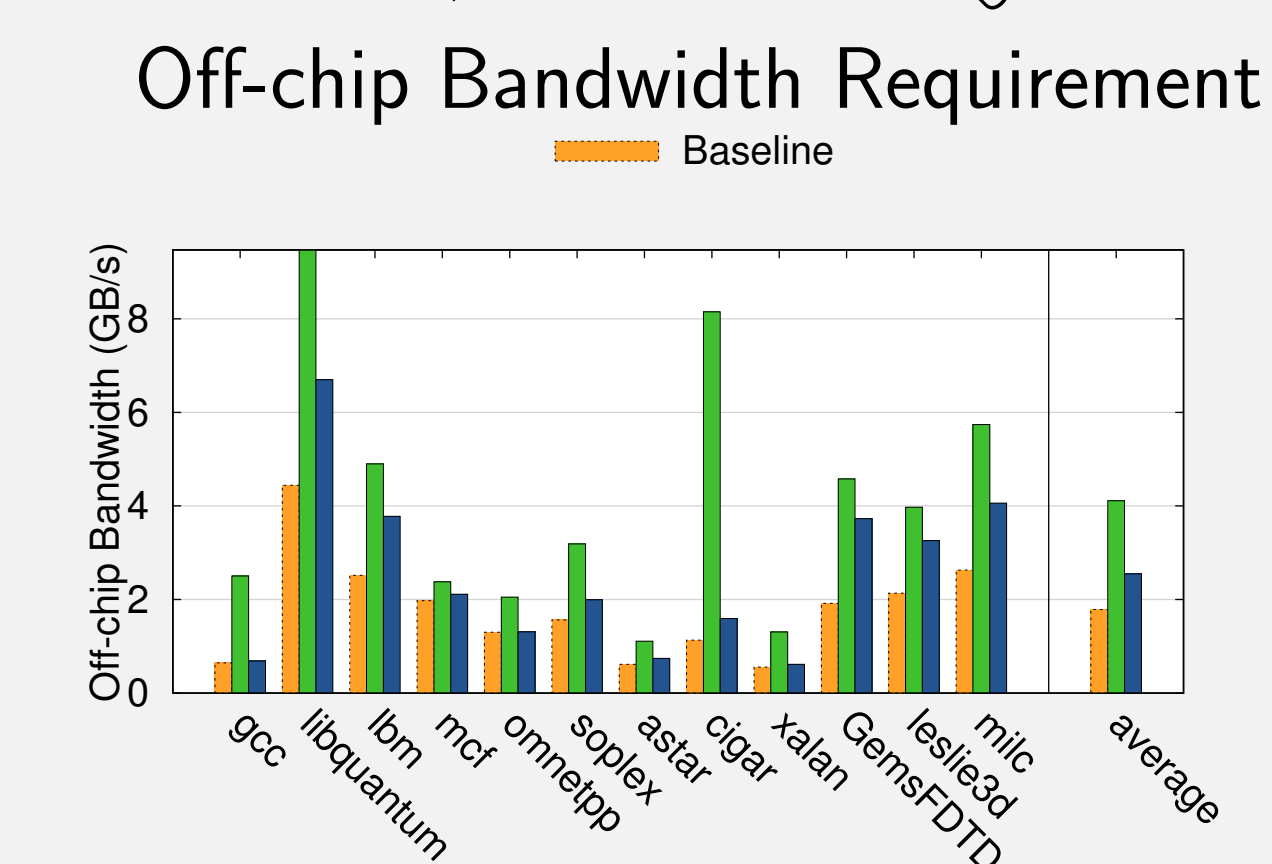
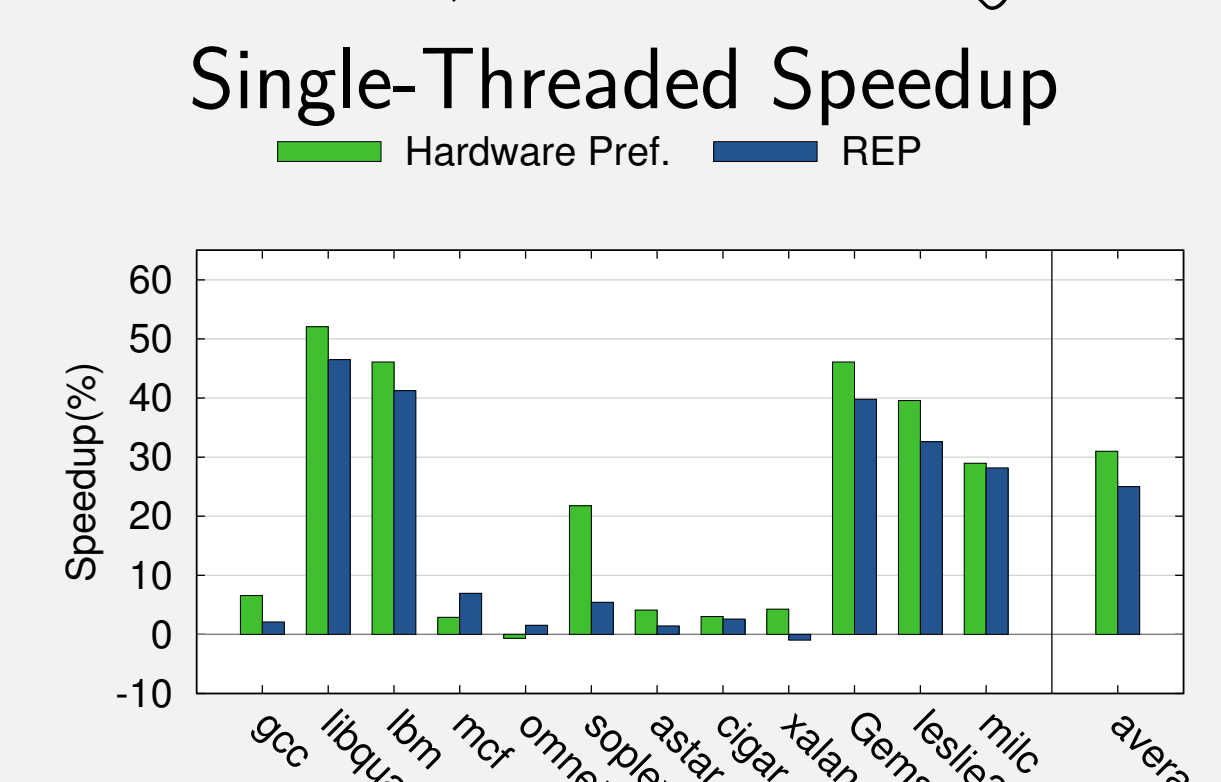
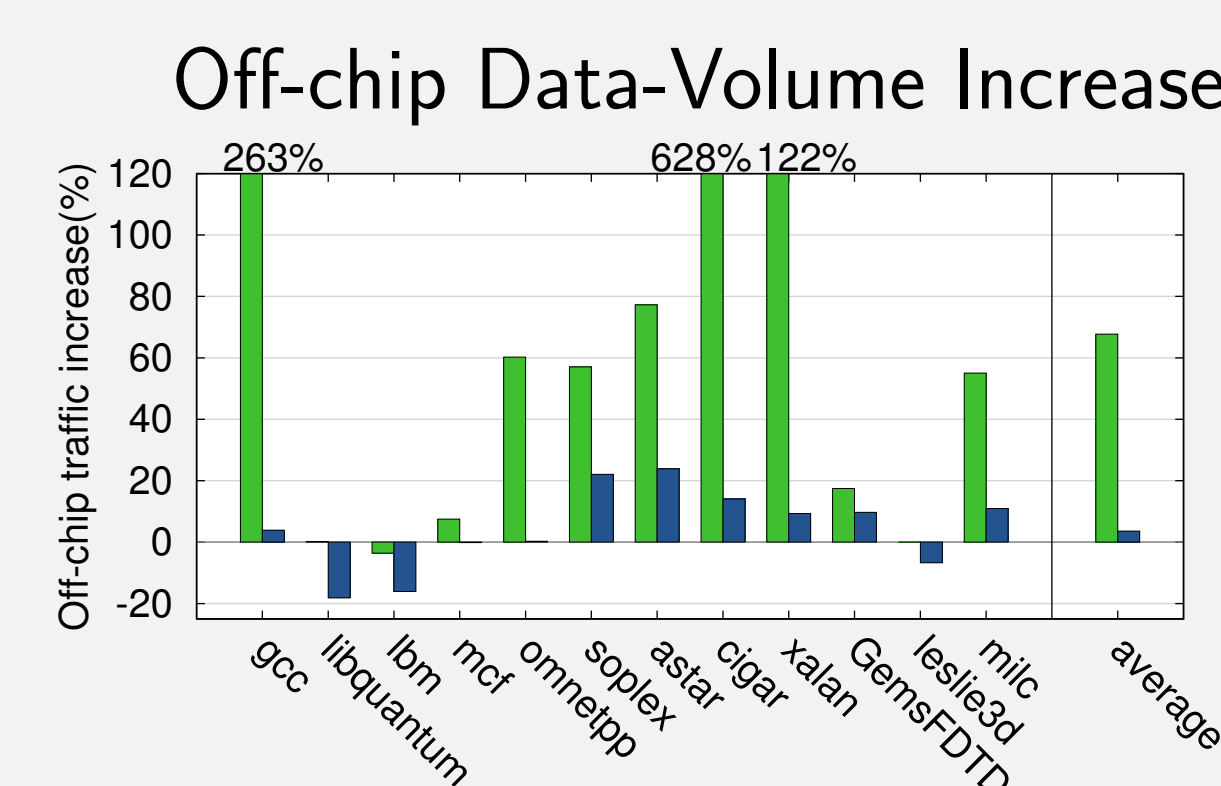
1- Insight: Hardware prefetchers require more shared resources for improving performance



2- Mixed Workload Throughput



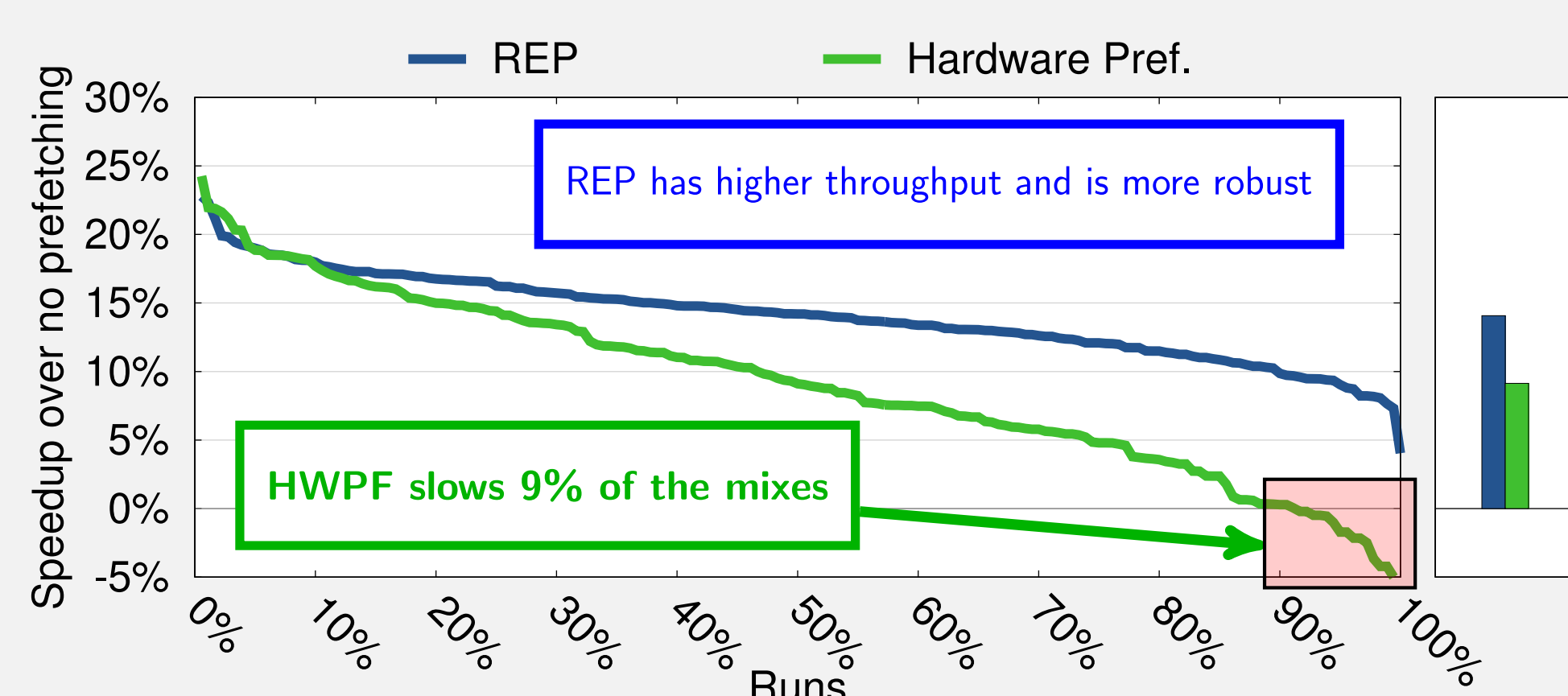
3- Single-Threaded Performance



- REP performance trails slightly behind hardware prefetching (within 5%)
- However, REP lowers off-chip traffic significantly, >60% on average
- Great potential for mixed workloads

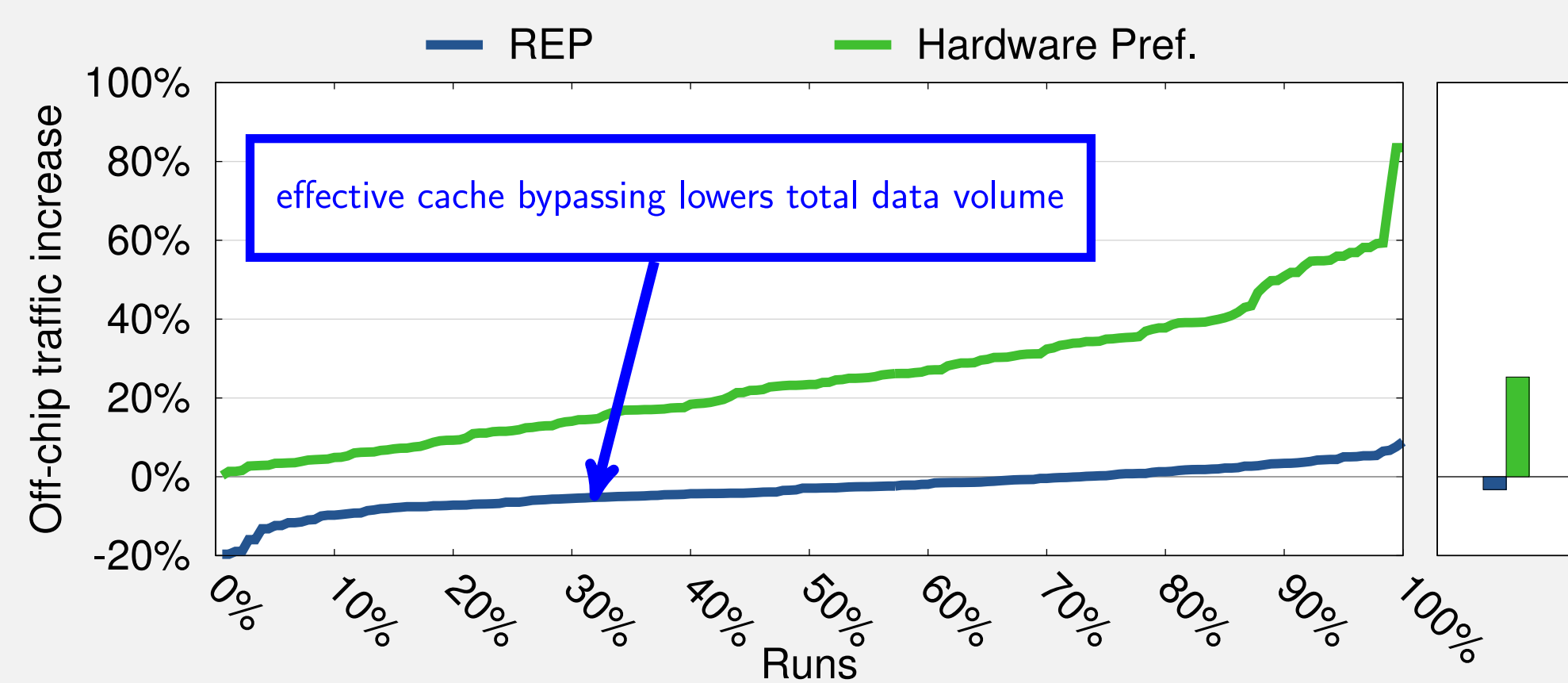
4- Mixed Workloads: Robust Prefetching Method

180 Mixed Workloads Throughput Performance (Intel)



Our scheme's performance is more robust than hardware prefetching across the 180 mixed workloads. REP performs 5% better on average. Cache bypassing helps lower data volume on average by 3% over the baseline.

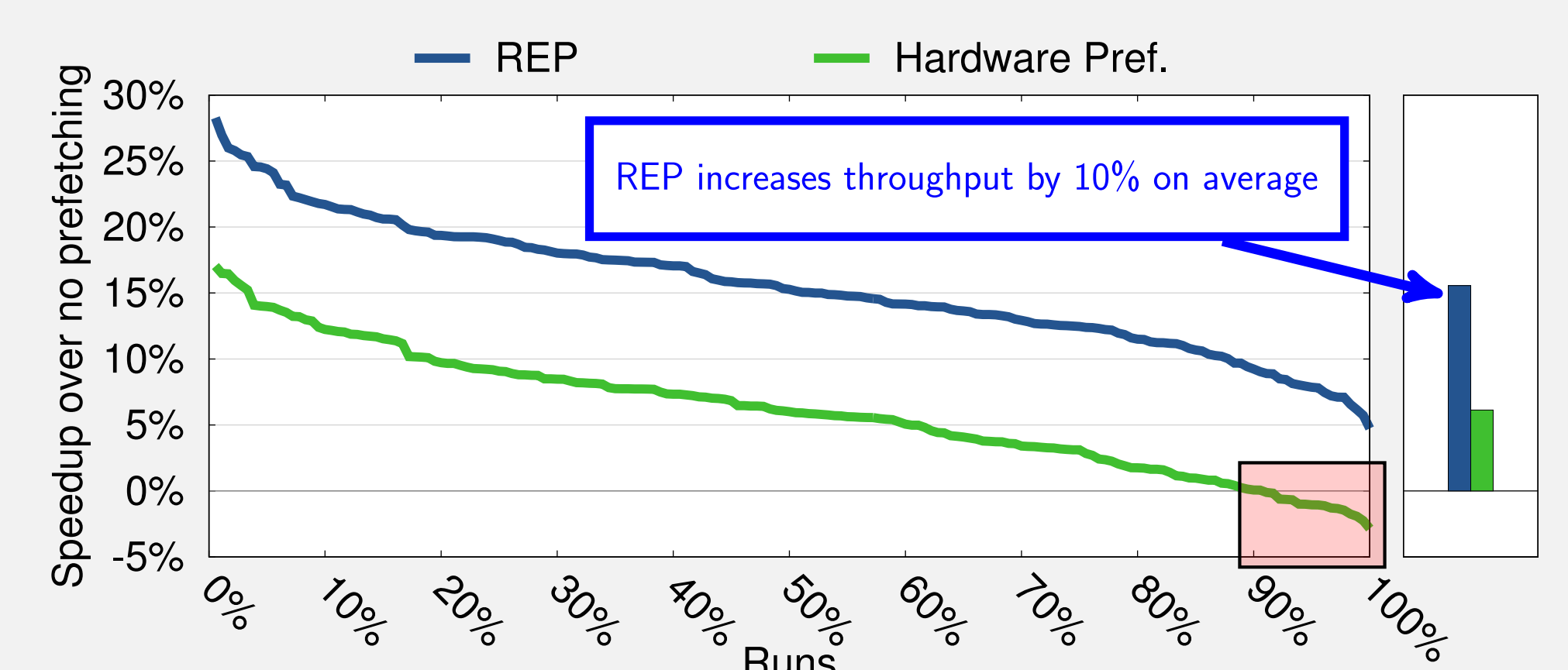
Increase in Total Data Volume



Intel i7-2600K (Sandybridge)

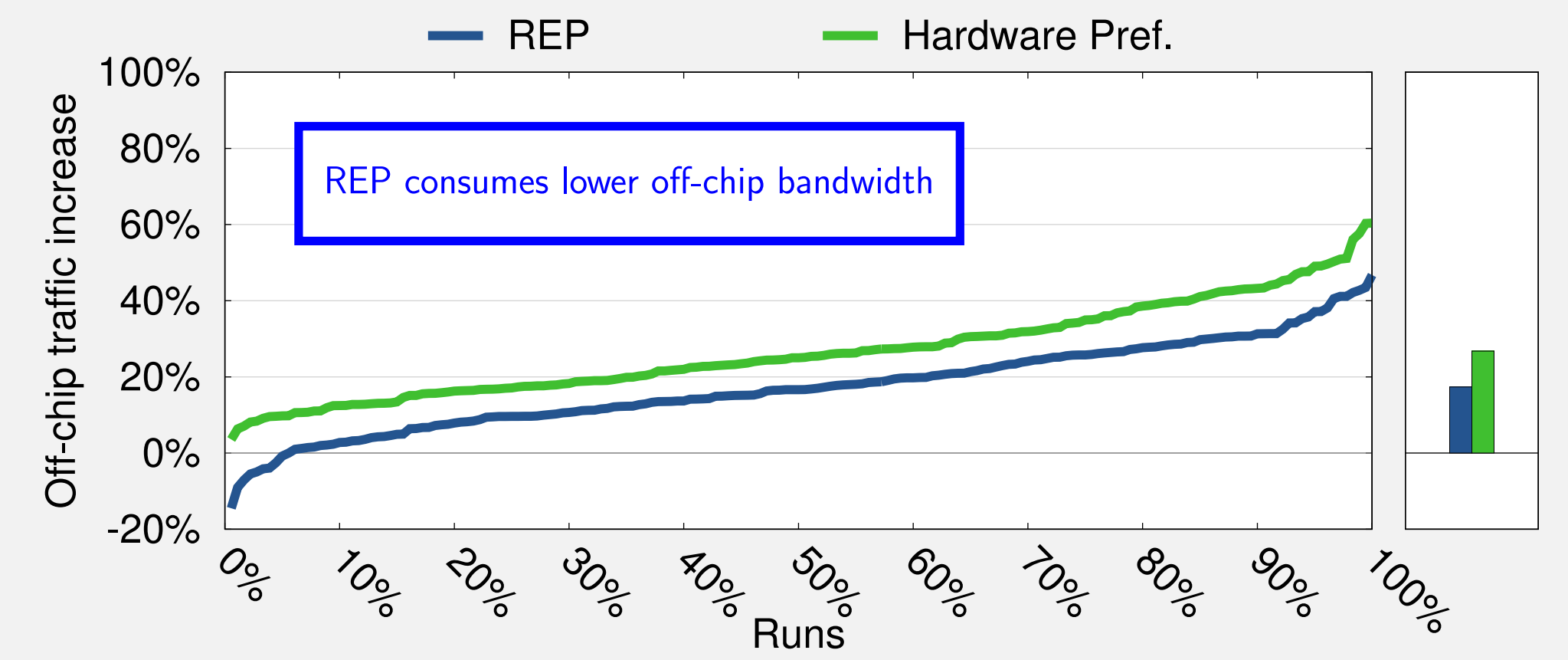
The benchmarks in Section 4 were selected to create 180 workload mixes. Each mix contains 4 different randomly chosen benchmarks that were run in parallel on the 4 cores of a Intel Sandybridge and AMD Phenom II processors. Different mixes stress shared resources differently, and help us explore how REP benefits performance under varying conditions. The graphs above compare REP's performance against hardware prefetching.

180 Mixed Workloads Throughput Performance (AMD)



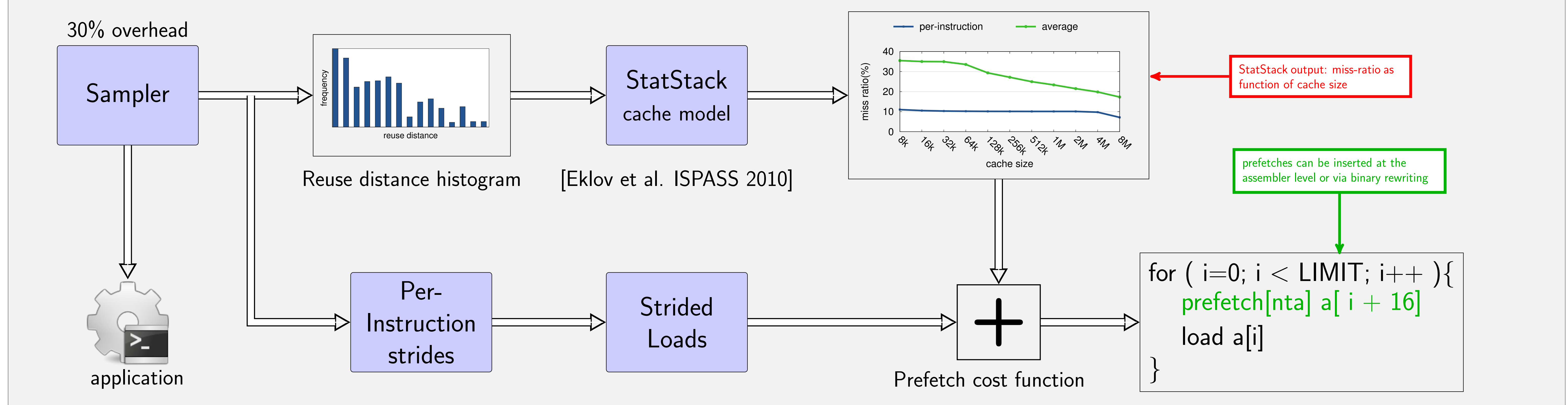
REP performs 10% better on average. On AMD, REP improved throughput performance across all mixed workloads, performing strictly better than hardware prefetching. REP consistently maintains less DRAM traffic.

Increase in Total Data Volume



AMD Phenom II X4 920

5- REP Framework: A fast software prefetching framework



6- Conclusions

This work investigates how a resource-efficient prefetching (REP) method can help improve throughput performance in multicores when shared resources are constrained. We propose an efficient method that 1) accurately prefetches the required data, 2) avoids (useless) speculative prefetching, and 3) employs cache bypassing to retain useful data in the higher level caches. In contrast to hardware prefetchers, REP is designed to maintain minimal off-chip traffic, and as a result avoids LLC pollution and lowers off-chip bandwidth demand.

This benefits throughput performance in multicores when several applications co-execute and share resources. Compared to state-of-the-art hardware prefetching on two high-performance commodity processors, REP performs up to 10% better on average. Our work highlights the importance of shared-resource friendly prefetching for optimizing multicore performance.