

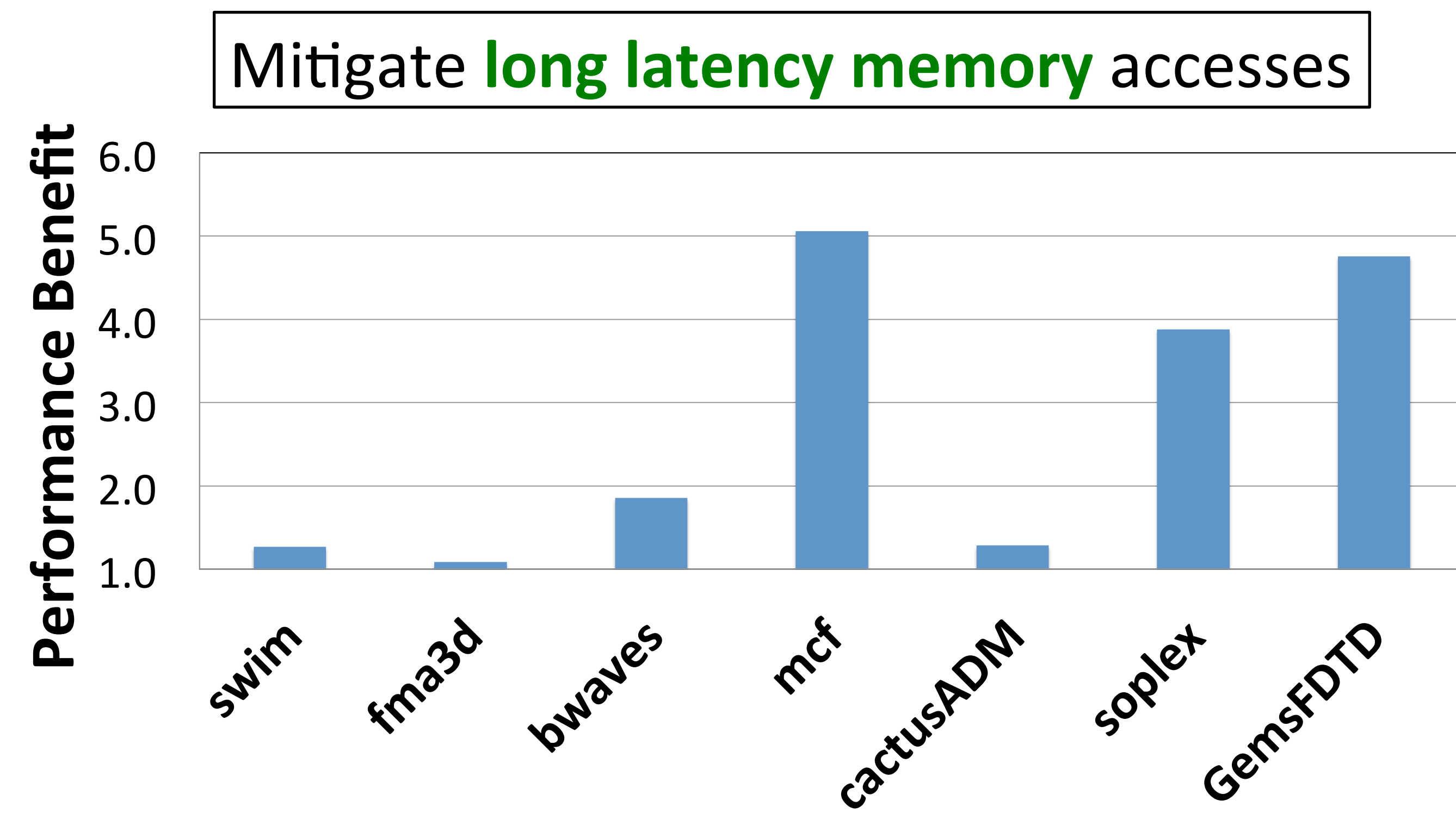


Rollback-Free Value Prediction with Approximate Loads

Bradley Thwaites, Gennady Pekhimenko, Amir Yazdanbakhsh, Jongse Park, Girish Mururu
Hadi Esmaeilzadeh, Onur Mutlu, Todd C. Mowry

Motivation

Perfect Prediction



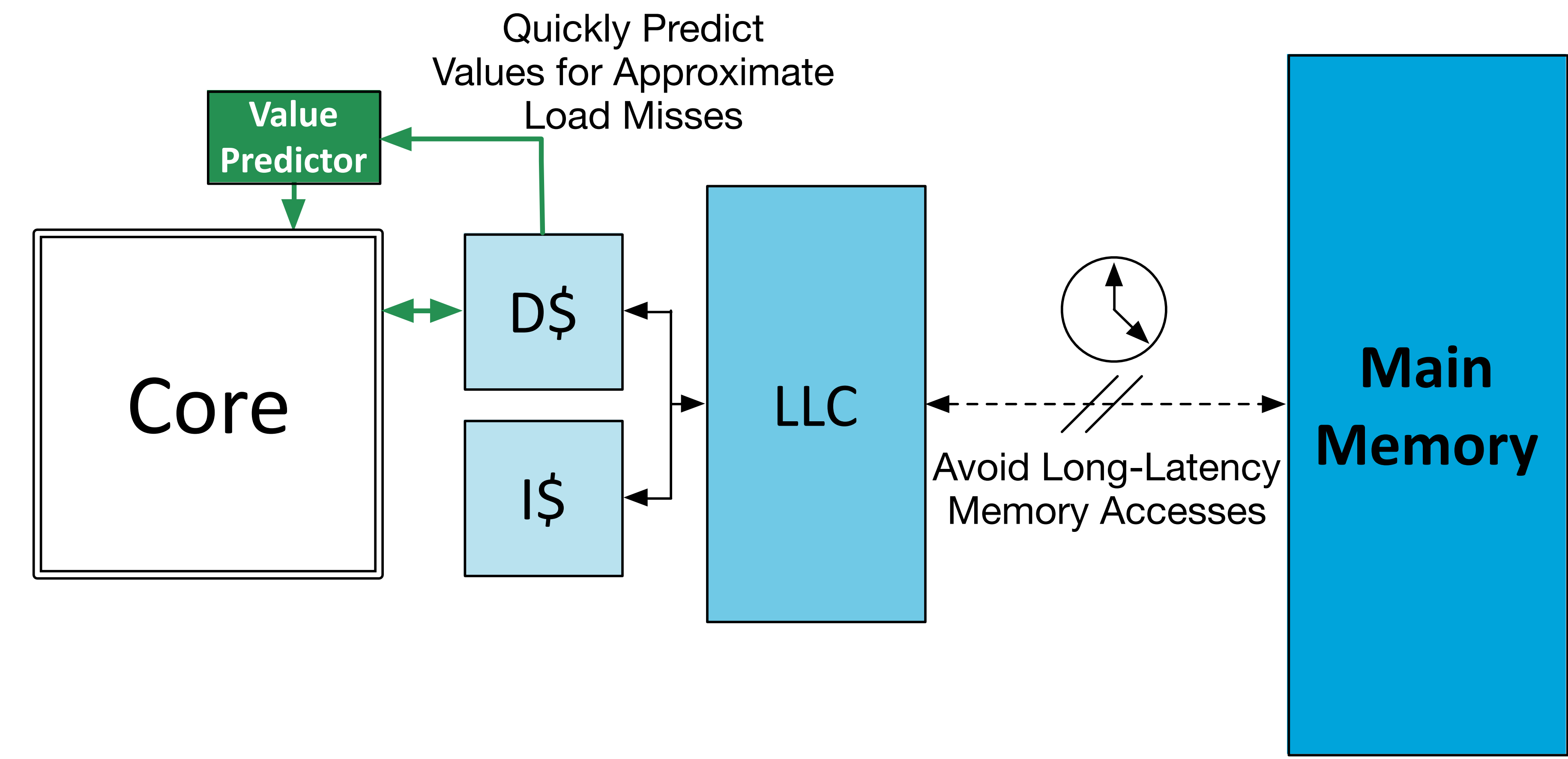
RFVP Overview

Microarchitecturally-triggered approximation

Predict the value of an approximate **load** when it **misses** in the cache

Do not **check** for mispredictions

Do not **rollback** from mispredictions



Design Principles

Maximize opportunities for performance and energy **benefits**

Minimize the adverse effects of approximation on **quality degradation**

Design Challenges

Target Performance-Critical Safe Loads

Profile-directed compilation

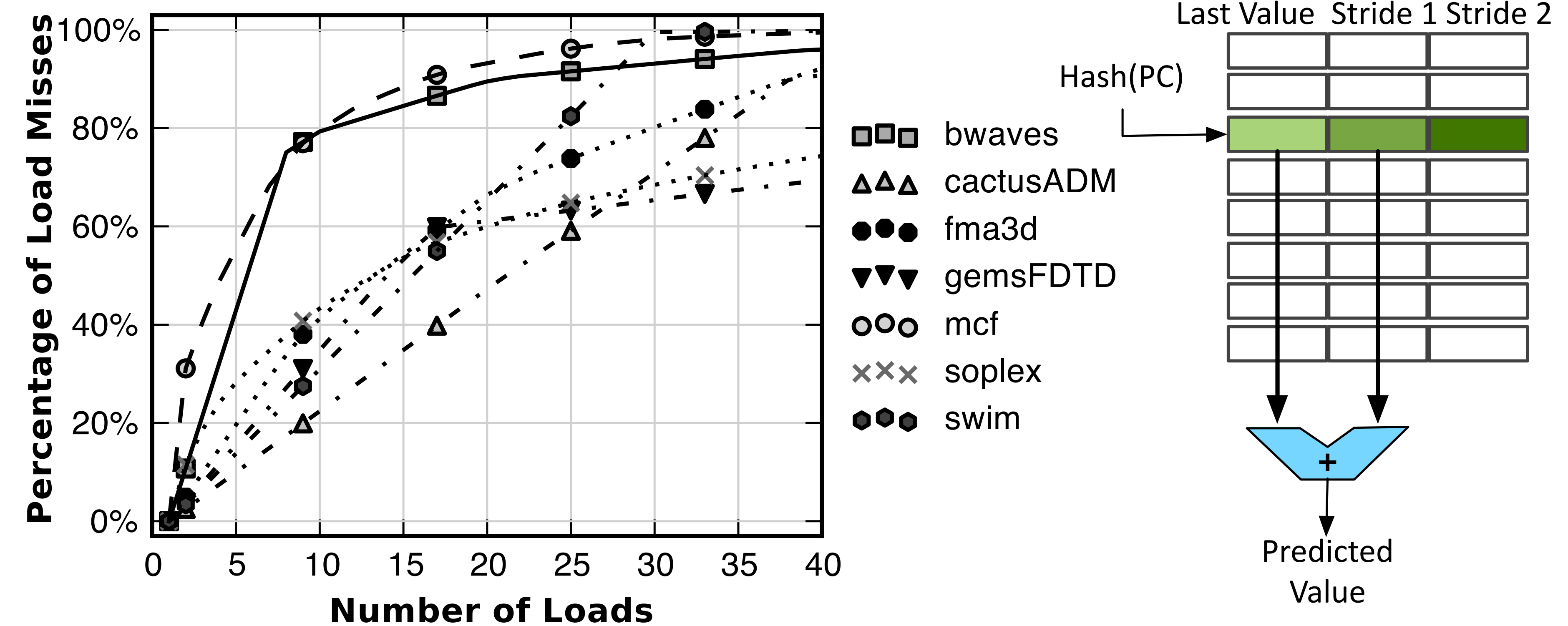
Usually, < 32 loads cause 80% of cache misses

Utilize Fast-Learning Predictors

Two-delta stride predictor

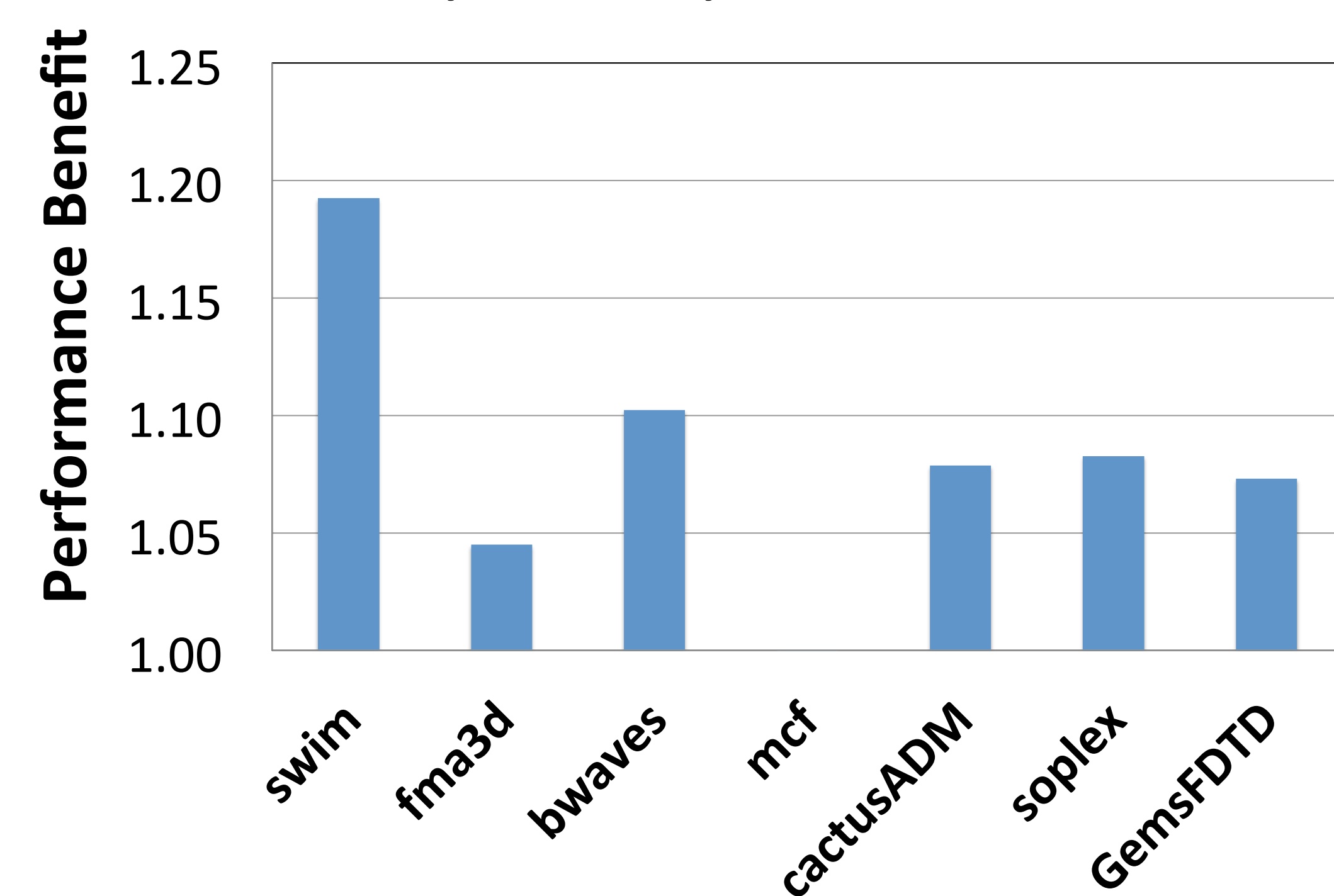
Prediction: table lookup plus an addition

Integrate RFVP with existing architecture

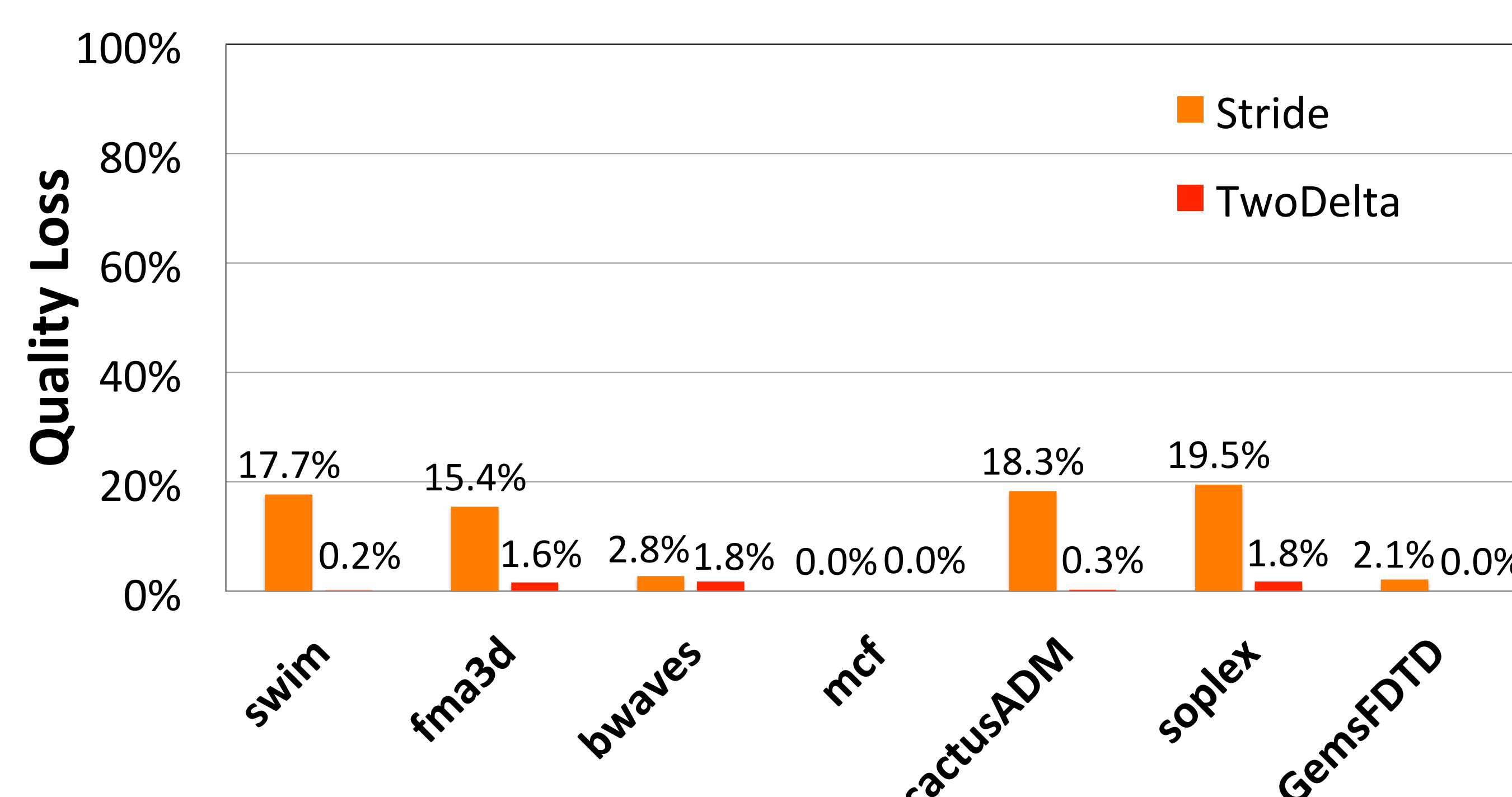


Key Experimental Results

2MB LLC, 4-Wide, Performance Results



Value Prediction - Quality Loss



Ongoing Work

Mitigate both **Memory Wall** and **Bandwidth Wall**

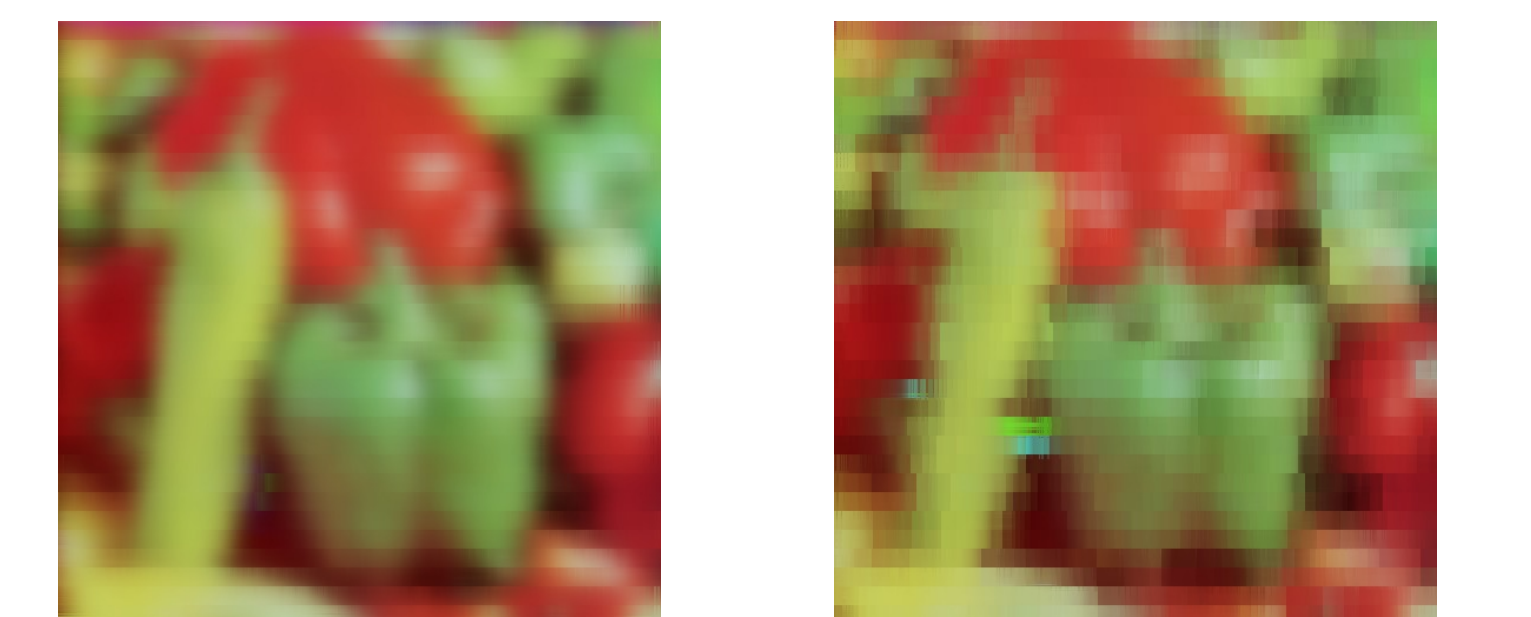
Extend rollback-free value prediction to GPUs

Drop a fraction of the missed requests

Preliminary results: Up to 2x improvement in

energy and performance

with only 10% quality degradation



(Performance Gain, Energy Gain, Quality Loss)