

Amphi: Practical and Intelligent Data Prefetching for the First-Level Cache

Xuan Tang*, Zicong Wang*, Shuiyi He, Dezun Dong[†], Xiangke Liao

College of Computer Science and Technology

National University of Defense Technology

Hunan, China

{tangxuan19, wangzicong, heshuiyi24, dong, xkliao}@nudt.edu.cn

Abstract—Data prefetchers play a crucial role in alleviating the memory wall by predicting future memory accesses. First-level cache prefetchers can observe all memory instructions but often rely on simpler strategies due to limited resources. While emerging machine learning-based approaches cover more memory access patterns, they typically require higher computational and storage resources and are usually deployed in the last-level cache. Other intelligent solutions for the first-level cache show only modest performance gains. To address this, we propose Amphi, the first practical and intelligent data prefetcher specifically designed for the first-level cache. Applying a binarized temporal convolutional network, Amphi significantly reduces storage overhead while maintaining performance comparable to the SOTA intelligent prefetcher. With a storage overhead of only 3.4 KB, Amphi requires only one-eighth of Pythia’s storage needs. Amphi paves the way for the broader adoption of intelligence-driven prefetching solutions.

Index Terms—data prefetcher, first-level cache, machine learning

I. INTRODUCTION

Hardware prefetching is a crucial technique for hiding long memory access latency. Among diverse levels of prefetching technologies, the data prefetching of the first-level data cache (L1D) presents substantial advantages on account of its proximity to the core and higher speed, thereby leading to significant performance.

Currently, there are two main categories of solutions for L1D prefetching: traditional methods and intelligent methods. Traditional prefetchers can easily capture regular and simple memory access patterns. These methods typically utilize tables to learn memory access patterns, which limits the learning capability due to constraints imposed by the size and representational ability of the tables. Intelligent methods often serve as auxiliary components in L1D to enhance performance. Some use perceptron to design a prefetch filter [4], and some create off-chip predictors to improve the overall performance [2]. Generally, intelligent prefetchers are often deployed in the last level cache (LLC). Although intelligent prefetchers can learn more complex or irregular memory access patterns compared to traditional prefetchers and achieve higher prediction accuracy, they are not deployed in L1D due to their significant storage overhead and resource consumption.

*Both authors contributed equally to this work.

[†]Corresponding author.

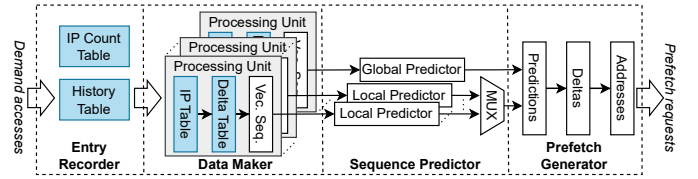


Fig. 1: Overview architecture of Amphi

We are committed to changing this situation by seeking an intelligent prefetcher that can be deployed in L1D. To further explore this, we have two key inspirations.

Inspiration 1: Data prefetching primarily revolves around the sequence of memory access patterns, which can be effectively optimized by sequence modeling analysis. Temporal convolutional network (TCN) is a typical sequence model, which can effectively capture long-range dependencies via dilated convolutions, and maintain stability in gradient propagation. The straightforward architecture and fewer hyperparameters also simplify design and tuning. Therefore, we choose TCN as the base model.

Inspiration 2: The introduction of binary neural network (BNN) has further advanced the field of model quantization, presenting new opportunities for intelligent prefetchers in L1D [5]. By quantizing the weights and activations of the model to binary, we can achieve substantial compression while maintaining performance.

To the best of our knowledge, **Amphi is introduced as the first practical and intelligent data prefetcher for the first-level cache.** Amphi integrates BNN with TCN to create a lightweight prefetcher, effectively balancing overhead and performance. We make the following contributions: 1) We propose Amphi, the first intelligent data prefetcher that can be deployed in the first-level cache. 2) Amphi is the first work to explore the application of BNN in data prefetching. 3) The evaluation demonstrates that Amphi outperforms the baseline by 23.79% in IPC with only 3.4 KB storage overhead.

II. METHODOLOGY

Figure 1 describes the overview architecture of Amphi, which mainly consists of four components: 1) Entry Recorder; 2) Data Maker; 3) Sequence Predictor; 4) Prefetch Generator.

Entry Recorder: The Entry Recorder consists of two tables: History Table (HT, 512 entries), and IP Count Table (IPCT, 16

