

Matryoshka: A Coalesced Delta Sequence Prefetcher

Shizhi Jiang

Yiwei Ci

Qiusong Yang

Mingshu Li





Presentation Structure

- Background
- Motivation
- Design & Implementation
- Evaluation
- Conclusion

Background

The "Memory Wall"



Year

Prefetcher

- learning access history
- predict future accesses
- Iow-overhead
- high-performance



Background

Delta Sequence

Delta = Current Addr - Last Addr

e.g. delta seq : 10, 5, 15

Recursive Lookahead



Advantages

- a) Shared Among Different Regions
- b) Record Access Order

Motivation

Why we need multi-matching ?

- Ideal Coverage
 - Always scattered in workloads
- Average Branch Number (Ideal Accuracy)
 Good in 3-delta or longer sequences

We can hardly decide the best delta sequence length!



Motivation

Better Form Stoing Patterns



Store 3 sequences separately

VLDP's Design



A Coalesced Delta Sequence

Our Design

It's obvious more efficient if we can find a method to coalesce variable-length sequences

Motivation

Storing Hot Patterns



Only a small part of deltas reappear frequently

Design

Reversed Coalesced Delta Sequences

The confidence (repeating times) of a sequence is vital for accurate prefetch! How to deal with them during coalescing?

- With same confidences
 - Discard confidences of short sequences
- With different confidences



Design

Dynamic Indexing Strategy & Overview



Delta Mapping Array Delta Sequence Sub-table

Dynamically record the mapping relations between deltas and table entries

Design

Consider both long and short matches

Adaptive Voting Strategy

$$S \ core_d = \sum_{i \in L} W_i \sum_{f \in M_i} Conf_i$$



In our configuration,

 W_i is the weight for each matched prefix of length *i*;

 $T_p = 0.5$ for prefetching into L1

 $W_2 = 3$ and $W_3 = 4$ for 2-delta and 3-delta matched prefixes.

L is the set of the lengths;

 M_i indicates the entries, containing the matched prefixes of length i, that generate d;

 $Conf_j$ is the confidence of the corresponding entry *j*;

D is the set of candidate deltas.

Implementation

Training

										L1 lo	ad acce	SS:
	History Table								PC:0x9a540			
							v			Addr:0x37caabd0170		
		PC tag Page tag			Last offset	t Last	Last delta sequence Valid					
		0x9a540 0x37caabd0		80		20,22,24		_1 ←				
	Cur delta seq: $24,22,20,26$ $(1) Cur offset=(0x37caabd0170>>3)&(2**9-1)=106$ $Cur delta=106-last offset=106-80=26$ $Cur delta seq= Reverse(I ast delta seq) \le 10 cur delta$									a		
		Pattern Table DSS							-			
			Vay 0 22 4 1		Delta Seq	Conf	Valid	Delta Seg	Conf	Valid		
		Way 0		4	1	$\xrightarrow{\text{Set 0}}$	18,8,36	3	1	25,20,36	1	1
(\mathfrak{I})		Way 1	16	5	1	$\xrightarrow{\text{Set 1}}$	20,18,22	2	1	20,18,24	1	1
24	⊢	Way 2	24	7->8	1	$\xrightarrow{\text{Set }2}$	22,20,26	7->8	1	0,0,0	0	0
		Way 3	26	8	1	Set 3	24, 2 2,28	8	1	24,19,24	3	1
						Fol	lowing seque	ence:	3)			
L <u>22,20,26</u>												

Implementation

Predicting



Evaluation: Methodology

Simulalor: ChampSim

Benchmark: 45 traces from SPEC 2017, CloudSuite

Table 1: Detailed Storage Overhead of Matryoshka

Structure	Entry	Field	Storage				
		PC tag (12 bits)					
	128 × 1	Page tag (8 bits)	7680 bits				
Listony Table		Last offset (9 bits)					
Flistory Table		Last delta sequence					
		(30 bits)					
		Valid (1 bit)					
Delta Manning		Delta (10 bits)					
Array	1×16	1×16 Confidence (6 bits)					
Allay		Valid (1 bit)					
Delta Seguence		Delta sequence (30 bits)					
Sub-table	16×8	Confidence (9 bits)	5120 bits				
Sub-table		Valid (1 bit)					
Candidate Array	128×1	Score (10 bits)	1280 bits				
Candidate Offset	32 × 1	Score (10 bits)	320 hite				
Array	52 ~ 1	5000 (10 0103)	520 0113				
$Total: 14,672bits \approx 1.79KB$							

Table 3: Prefetcher Overhead

VLDP	SPP+PPF	Pangloss	IPCP	Matryoshka
48.34 KB	48.39KB	45.25KB	740 B	1.79KB

we expand VLDP' s storage capacity to 48 KB

Evaluation: Single-core

Matryoshka yields the best geometric mean speedup of 53.1% over the non-prefetching system, an improvement of 6.5% over IPCP.

- Compared with the other prefetchers, Matryoshka achieves the highest average coverage (57.4%) at L1, which exceeds the second best IPCP by 6.0%.
- The average overprediction rate for Matryoshka is 20.6%, which is the lowest.



Evaluation: Multicore

In the 4-core system, Matryoshka offers a 32.2% speedup over the baseline





■ IPCP ■ VLDP ■ Pangloss ■ SPP+PPF ■ Matryoshka

Evaluation: Sensitivity Study



Compared with the other four prefetchers, Matryoshka still obtains better performance in different configurations of the cache system.

Conclusion

- Coalesced Delta Sequence
 - Simpler Structure for Multiple Matching
- Adaptive Voting Strategy
 - Ensure both coverage and accuracy
- Dynamic Indexing Strategy
 - Keep key information