Deep reinforcement learning in loop fusion problem

Group 15: Siqi Shao, Yixin Shi, Ying Yang 04/03/2023

https://doi.org/10.1016/j.neucom.2022.01.032



Problem Statement: <u>Loop Fusion</u> is a loop transformation in which several loops are combined to form a loop.

Goal:

- Reduce overhead
- Increase parallelism
- Enhance data reuse

for i in range(100):	
A[i] = B[i]+1	
	Synchronization
for i in range(100):	
A[i]*=2	

Fig1. A synchronization requisite between loops in a parallel system



Problem Statement: <u>Loop Fusion</u> is a loop transformation in which several loops are combined to form a loop.

Goal:

- Reduce overhead
- Increase parallelism
- Enhance data reuse

for i in range(100):	
A[i] = B[i]+1	Curahaniatian
fan i in nanna (100):	Synchronization
for I in range(100):	
A[I]^=2	

Fig1. A synchronization requisite between loops in a parallel system

Important technique for optimizing the performance of scientific or computational algorithms that involve repetitive operations.

Loop Fusion is a NP-hard problem

Factors to consider in loop fusion.

- Data dependencies
- Data reuse
- Loops' type
- Computer system's register size



Kelated code	
for i in range(100):	
A[i] = i	
for i in range(100):	
B[i]=A[i]+2	
F[i]=C[i]	
for i in range(100):	
C[i]=4	
H[i]=i*3	
for i in range(100):	
H[i]=G[i]	
for i in range(100):	
K[i]=A[i-1]+G[i]	
Loops	
Dependence edges	

Loop 1 and Loop 5: Loop carrier dependence based on array A.

Related Work

Year Paper's g	Paper's goal	Number	mber Graph's	Loop order	Modeling method	Improvement factors		
	C	of types node			parallelism	data reuse	Register size	
1984	Improving parallelism	2	Loop and instruction	Topological order of loop independent dependence graph	hierarchical dependence graph (HDG)	1	1	×
1993	Provide two distinct methods for improving parallelism and data	2	loop	Breath-first-search with the parallel type's priority	Directed acyclic graph	1	1	×
2004	Improving data reuse considering register size	2	loop	Topological order with the priority of removing the edges with maximum weight	Directed weighted acyclic graph	×	1	1
2020	Improving data reuse with maintaining parallelism	2	permutation	Topological order considering a coloring heuristic	Fusion conflict graph (FCG)	1	1	×
2020	Improving parallelism and data reuse considering register size	>2	loop	The sequence generated by evolutionary algorithm	Directed acyclic graph	1	1	1

Hard to achieve all three improvement factors; Long time to solve the problem.

Deep Reinforcement Learning Loop Fusion

- Find a proper order for loop fusion
- Enhance parallelism
- Consider the system's register size as a boundary for fusion
- First use deep reinforcement learning to solve loop fusion in a short span of time.

Proposed Algorithm

- Goal: determine the best order of loops for merging
 - the runtime of the resulting loops is optimal
- Main framework

Deep reinforcement learning:

Repeatedly:

- Encode the input.
- Decode and Estimate the fusion sequence of loops.
- Evaluate the resulted sequence based on the parallelism improvement of fused loops according to the system register size.

Fig. 3. The framework of the proposed model.

Output:

Predicted fusion order

Preprocessed input:

Includes information as:

- Type of loops
- Number of arrays of that loop
- Data dependencies between them
- Number of arrays that can be stored in the computer system register

- Draw a loop dependency graph
 - Sequential loops vs Parallel loops
 - Whether it can execute multiple iterations simultaneously
 - Dependency edges
 - Fusion preventing edges



- Create an adjacency matrix
 - n × n matrix
 - (i, j) -> loop i to loop j
 - hold floating-point numbers
 - Include features
 - Type of loops
 - Data dependencies between loops
 - Number of arrays of loops
 - Number of arrays that can be stored in the computer system register





- Create an adjacency matrix
 - Main diagonal: information about the type of each loop
 - Integer part: the size of the system register





- Create an adjacency matrix
 - Main diagonal: information about the type of each loop
 - Fractional part: number to the loop's type



The system register size

- Create an adjacency matrix
 - The other matrix elements: dependencies between the loops
 - Integer part: type of dependence



Required

register size

for

executing

i loop

points

Dependence

type

between

i and j loops

- Create an adjacency matrix
 - The other matrix elements: dependencies between the loops
 - Fractional part : the number of register units required to execute the loop i

Required

register size

for

executing

i loop

points

Dependence

type

between

i and j loops



Encoder

- Input: loop specifications generated by input preprocessing
- Output: extract a set of vectors, each of which specifies a loop
 - Characteristics of each loop
 - Correlation between loops
- Produce the action space of the decoder



Encoder

- Sublayer: multi-head attention
 - Attention Ο



- Multi-head Ο
 - Attend to different parts of the input matrix "V" simultaneously
 - Use multiple sets of query, key, and value vectors.
 - Increase model capacity and performance

 $head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$

Attention weight

 $MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)$

Decoder

- Input
 - Embedded loops vectors generated by encoder
- Output
 - The probability distribution of order for fusing loops



Decoder

- Calculate query vector
 - State representation in every single step
 - \circ ~ Consider the last three added loops (actions) at time t

$$q_t = ReLU(W_1a_{\pi(t-1)} + W_2a_{\pi(t-2)} + W_3a_{\pi(t-3)}) \in \mathbb{R}^d$$

- Calculate the probability distribution on the remaining loops
 - O pointing mechanism

$$\forall i \leq n.u_i^t = \begin{cases} v^T \tanh(w_{ref}a_i + w_qq_t) & \text{if } i \notin \pi(0) \dots \pi(t-1) \\ -\infty & \text{otherwise.} \end{cases}$$

 $p_{\theta}(\pi(t)|\pi(< t).s) = softmax(C \tanh(u^t/T))$

ai: added loops (actions)qt: query (state representation)w: two attention matrixv: attention vectorT: temperature hyper-parameter



1

Reward function

• Minimize

 $r(\pi|s) =$ number of topological order violations in the considered sequence + number of consecutive loops in the considered sequence that are not fusiable + number of times that fusible loops' required register size passes the system's register size + One tenth of the total weight of the edges between nodes in the sequence

- Encourage the agent to select loop orders that
 - maintain topological order
 - have more fusible paris
 - don't exceed the available system register size
 - decrease the dependencies

Reward function

 $r(\pi|s) =$ number of topological order violations in the considered sequence + number of consecutive loops in the considered sequence that are not fusiable + number of times that fusible loops' required register size passes the system's register size + One tenth of the total weight of the edges between nodes in the sequence



Experiment

- Trained with a set of 1000 test cases, each with a maximum of 8 loops owing the size of the available benchmarks
- 8 * 8 matrices as input
- Compare with previous methods from two perspectives:

1. Output in terms of fusion sequence and the resulting **number of loops**

2. Execution time in sequential and parallel

Result

Table 4

Comparing the proposed method with comparative works in terms of fusion sequence and the resulted loops.

[Ref.]	Original loops sequence	Referred paper's result		DRLLF's result
[31]	12345678	1-2 3-4 5-6 7-8		1-2 3-4 5-6 7-8
	Number of loops: 8	Number of loops: 4		Number of loops: 4
[32]	12345678	1-3 2-4-6 5-8 7		1-3 2-4-6 5-8 7
Contraction of the second seco	Number of loops: 8	Number of loops: 4		Number of loops: 4
[33]	12345678	1-2-8 3-4-5-6 7	Fewer loops	1-2-3-87-4-5-6
	Number of loops: 8	Number of loops: 3		Number of loops: 2
[8]	123456	21-34-56		21-34-56
	Number of loops: 6	Number of loops: 4		Number of loops: 4
[35]	12345	123-4-5		1-23-4-5
- Forest	Number of loops: 5	Number of loops: 3		Number of loops: 2
[36]	1234567	1-2-3-4-675		1-2 7 3-4-5-6
	Number of loops: 7	Number of loops: 3		Number of loops: 3
[41]	123456	1-2-4 3-5-6		1-2-4 3-5-6
	Number of loops: 6	Number of loops: 2		Number of loops: 2
Ben6[41]	1234567	1-3 4-7 2-6 5		1-3 4-7 2-6 5
	Number of loops: 7	Number of loops: 4		Number of loops: 4
LL8[41]	123456	1-2-3-4-5-6	more loops	1-2-3 4-5-6
	Number of loops: 6	Number of loops: 1		Number of loops: 2
LL18[41]	123456	1-2 4-3-5-6		1-2 4-3-5-6
	Number of loops: 6	Number of loops: 2		Number of loops: 2

Result



Pros & Cons

- Pros
 - Effective deep reinforcement learning
 - Consider the system's register size
 - Enhance parallelism
- Cons
 - Long training time (3 weeks for training!!)
 - May not generalize well to new datasets
 - Should consider more complex real-world dependencies

Future Work

- Explore the use of other deep learning techniques in loop fusion optimization, such as **CNN or RNN**.
- Investigate the application of the proposed method to other loop optimization problems, such as **loop tiling or unrolling**.
- Exploring the **use of unsupervised learning** methods to reduce the reliance on labeled training data and improve the scalability of the proposed approach.
- Investigate how this approach can be **integrated into existing parallelizing compilers** and evaluate its performance on a wider range of benchmarks.

Thanks

Q & A