

Efficient and Portable ALS Matrix Factorization for Recommender Systems

Jing Chen¹, Jianbin Fang¹, Weifeng Liu², Tao Tang¹, Xuhao Chen¹, Canqun Yang¹

*College of Computer, National University of Defense Technology*¹

*Niels Bohr Institute, University of Copenhagen*²

Contents

- Background
- Motivations
- Design and Implementation
- Experimental Setup
- Performance Results
- Conclusion

· 1.BACKGROUND ·

1.1 Recommender systems

System Goal :

build a model

train with observed incomplete rating data;

and predict preference over items not rated.

Recommendation Approaches :

Matrix factorization (MF), nearest-neighbor ...

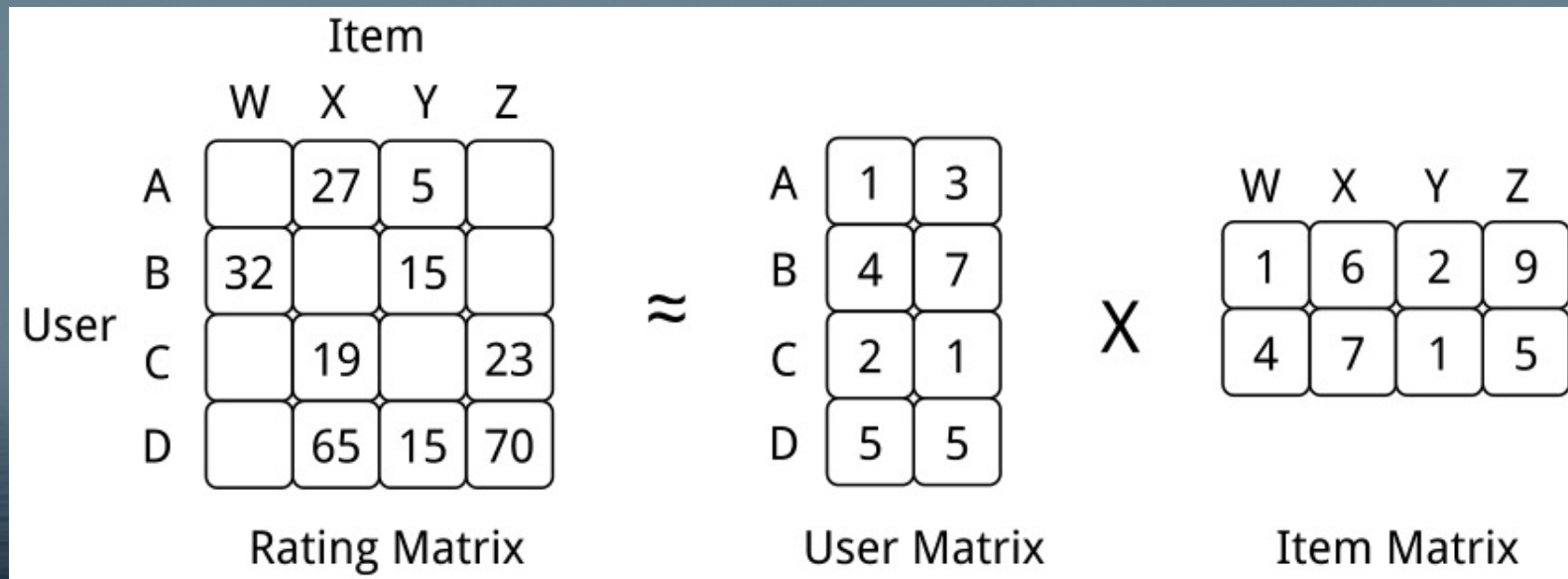
Popular Algorithms of Matrix Factorization :

ALS (Alternating least squares), SGD (Stochastic gradient descent),

CCD (Cyclic coordinate descent) ...

1.2 Matrix Factorization

- Input : Rating matrix between users and items, $R(m \times n)$
- Output : $X(m \times k)$ matrix and $Y(n \times k)$ matrix so that $r_{ui} \approx x_u y_i^T$



1.2 Matrix Factorization

- Input : Relation matrix between users and items, $R(m \times n)$
- Output : $X(m \times k)$ matrix and $Y(n \times k)$ matrix so that $r_{ui} \approx x_u y_i^T$
- minimize the regularized squared error to obtain X, Y

$$L(X, Y) = \sum_{u, i \in \Omega} (r_{ui} - x_u^T y_i)^2 + \lambda (|x_u|^2 + |y_i|^2)$$

x_u^T : the u^{th} row vector of matrix X

y_i : the i^{th} column vector of matrix Y

Ω : all the nonzero ratings of matrix R

λ : regularized coefficient (to avoid over-fitting)

1.3 ALS

Principle : to keep one fixed while calculating the other

1. We minimize the equation over X while fixing Y , the function becomes,

$$L(X) = \sum_{i \in \Omega_u} (r_{ui} - x_u^T y_i)^2 + \lambda |x_u|^2$$

2. Calculating derivative of x_u and let the partial derivative equal zero,

$$x_u = (Y^T Y + \lambda I)^{-1} Y^T r_u$$

3. In a same way, $y_i = (X^T X + \lambda I)^{-1} X^T r_i$

4. ALS iterates until it reaches the maximum specified cycles or error rate.

2.MOTIVATIONS

2.1 Motivations

Observation 1 :

ALS on CPUs runs faster than on GPUs.

- ✓ 11X faster on the CPU than on the GPU.
- ✓ Restructure the algorithm
- ✓ Customize optimizations according to the architectural specifics.

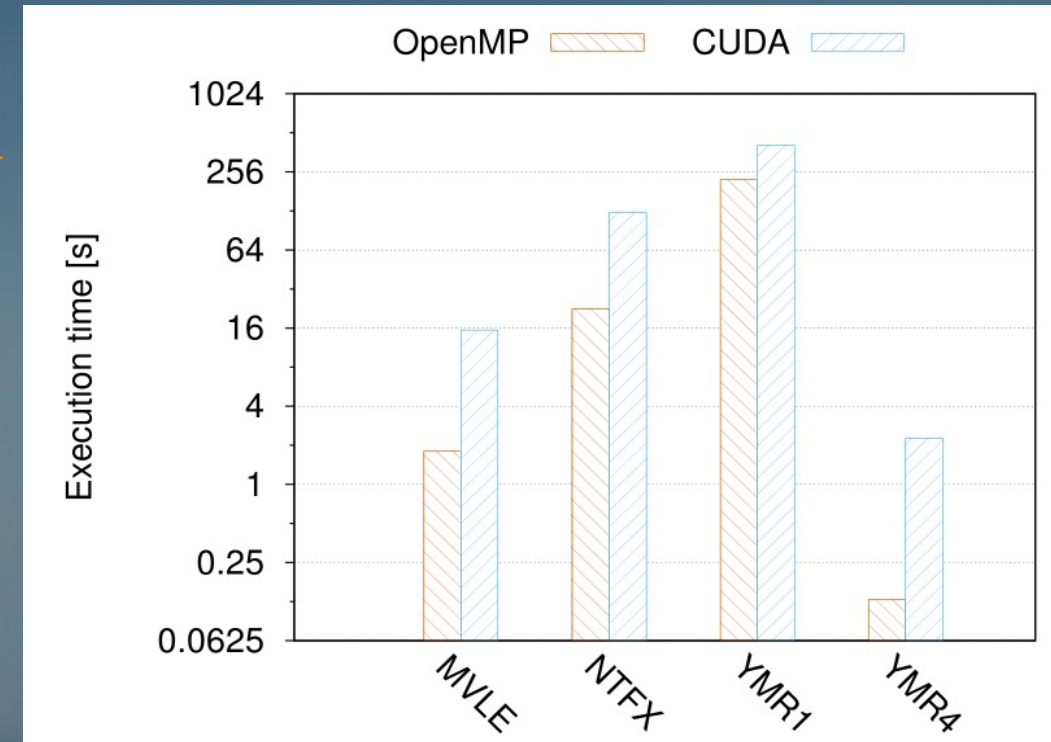


Fig 1. Performance comparison of an OpenMP implementation on a 16-core CPU versus a CUDA implementation on K20C.

2.2 Motivations

Observation 2 :

The current implementation cannot run on the coprocessors, such as Intel Xeon Phi.

Various Platforms : GPUs, MICs, FPGAs, DSPs ...

The current implementations cannot be offloaded to run on FPGAs.

Porting is time-consuming and error-prone.

- ✓ Speed
- ✓ Portability

3. DESIGN and IMPLEMENTATION



Thread Batching Parallelization

Baseline Design :

using **one thread** to update a row of X or a column of Y.

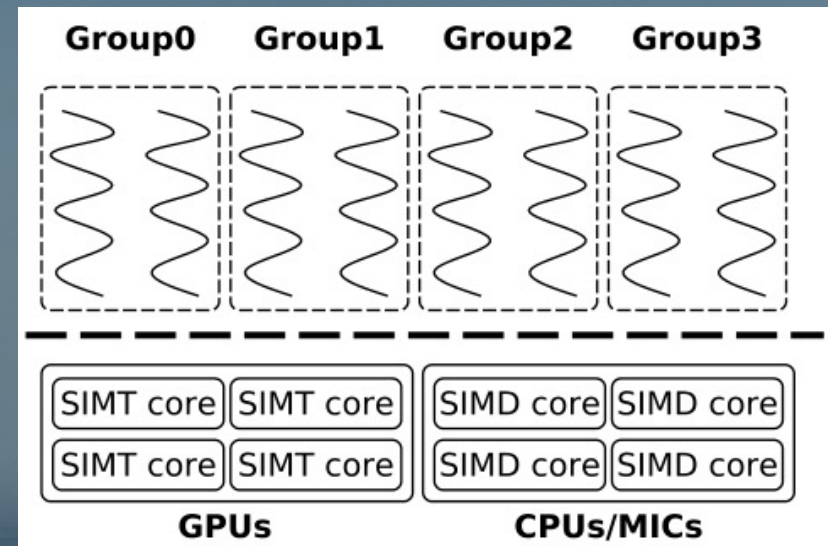
unaware of the hierarchical thread organization in CPU / GPU / MIC

Problem :

- unbalanced thread use
- random memory access

Solution :

- using thread batching technique
- wrap a branch of threads to deal with a row / column



Architecture-specific Optimizations

■ Using Registers

Modern GPUs feature abundant registers (small accessing latency).

Tesla K20C : 256KB registers in each SM.

```
1 float sum[k*k]={0};
2 for (int i = lx; i < k; i+=ws){
3     for (int j = i; j < k; j++) {
4         for (int z = 0; z < omegaSize; z++){
5             int d = col_idx[row_ptr + z] * k;
6             sum[i*k+j] += Y[d + i] * Y[d + j];
7         }
8         smat[(j*k)+i] = sum[i*k+j];
9         smat[(i*k)+j] = sum[i*k+j];
10    }
11 }
```

(a) Original code.

```
1 float sum0=0,sum1=0,sum2=0,sum3=0,sum4=0;
2 for (int z = 0; z < omegaSize; z++){
3     int d = col_idx[row_ptr + z] * k;
4     if(0<=lx<k) sum0 += Y[d + lx] * Y[d + 0];
5     if(1<=lx<k) sum1 += Y[d + lx] * Y[d + 1];
6     if(2<=lx<k) sum2 += Y[d + lx] * Y[d + 2];
7     if(3<=lx<k) sum3 += Y[d + lx] * Y[d + 3];
8     if(4<=lx<k) sum4 += Y[d + lx] * Y[d + 4];
9     ...
10 }
11 // updating the smat matrix
```

(b) Unrolling the code.

Figure 2. An example of unrolling the code to calculate $Y^T Y$, where $k=5$.

Original version: private array $\text{sum}[k*k]$ for each thread

Unrolling version: **k registers** for each thread block are enough

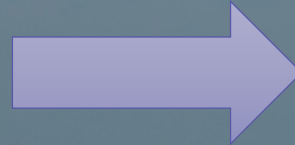
Architecture-specific Optimizations

■ Using Scratch-pad memory (Local memory in OpenCL)

- is a high-speed memory unit located on-chip
- *data sharing* in a same thread block
- increase *data moving bandwidth* between off-chip and on-chip memory

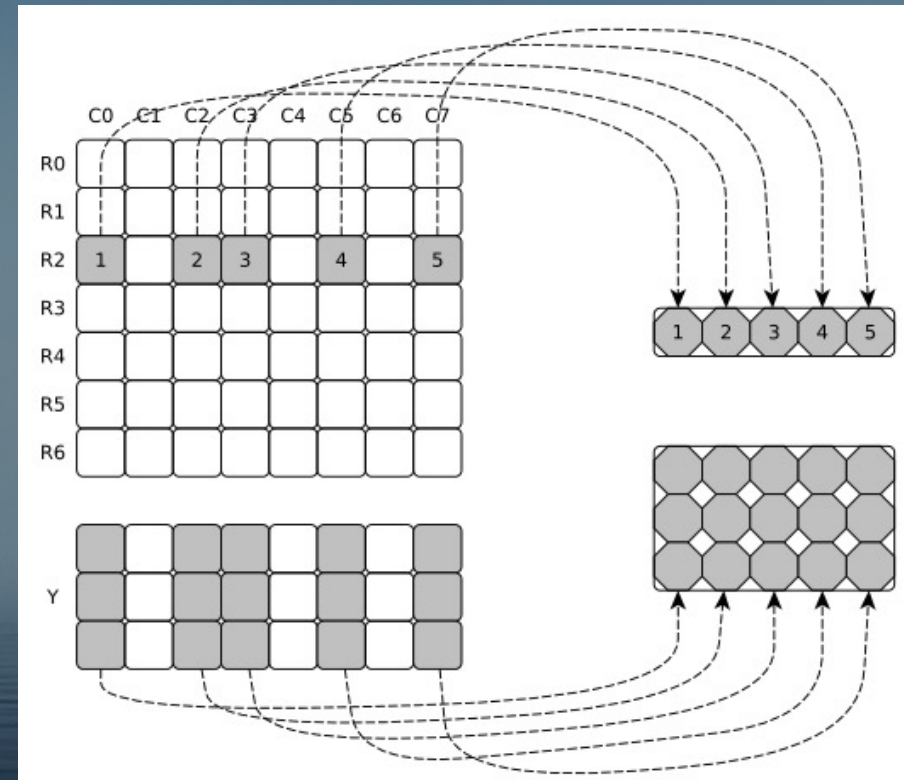
Sparsity of R matrix

Incontiguous Data

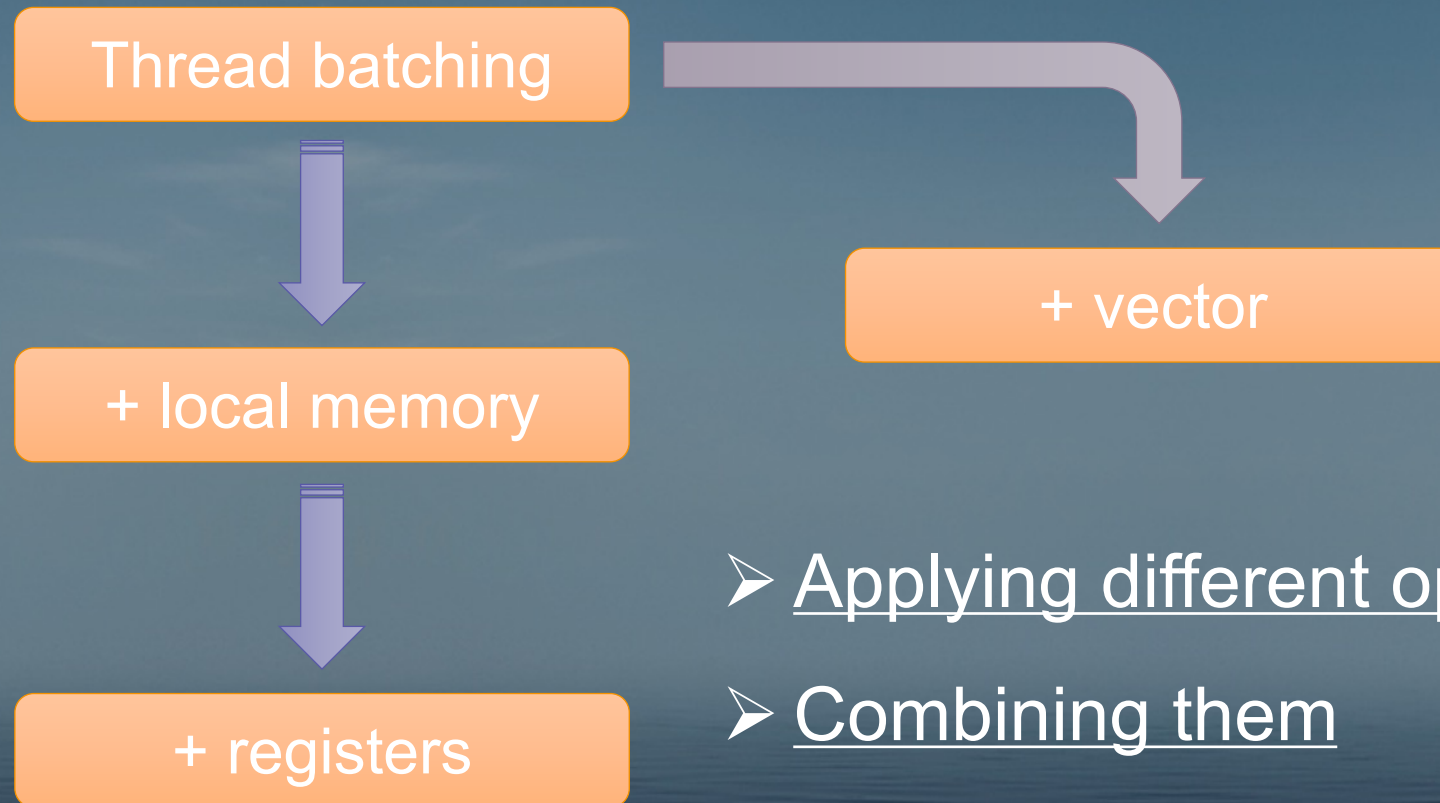


■ Using Vector Units

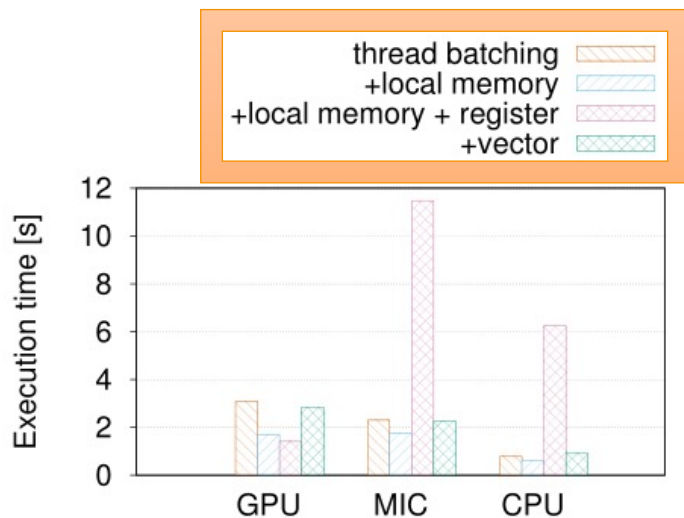
- CPU
- MIC



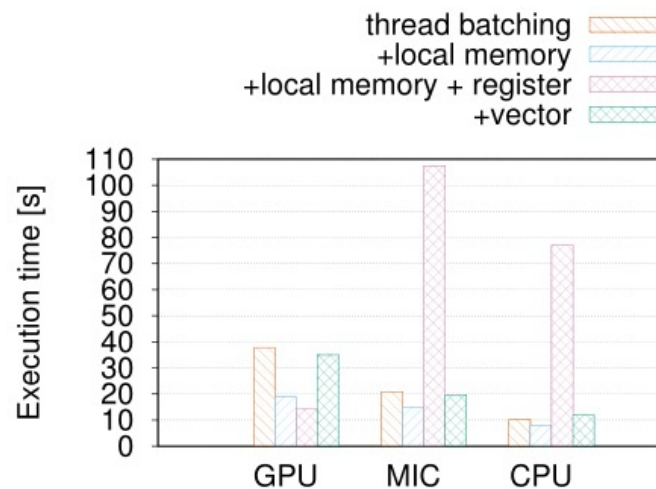
Code Variant Selection



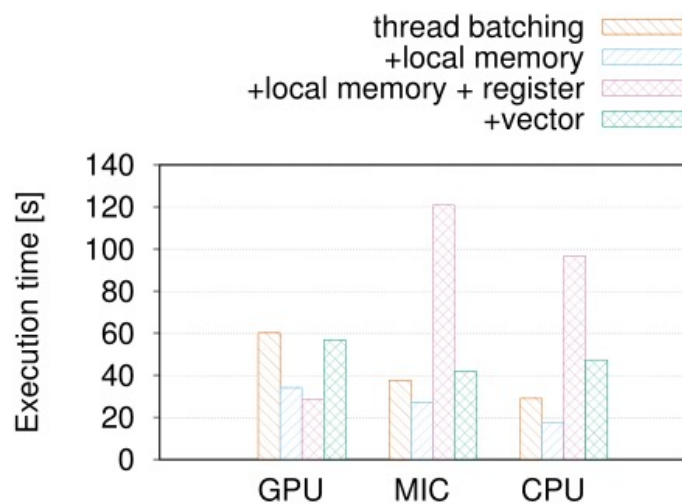
- Applying different optimizations
- Combining them



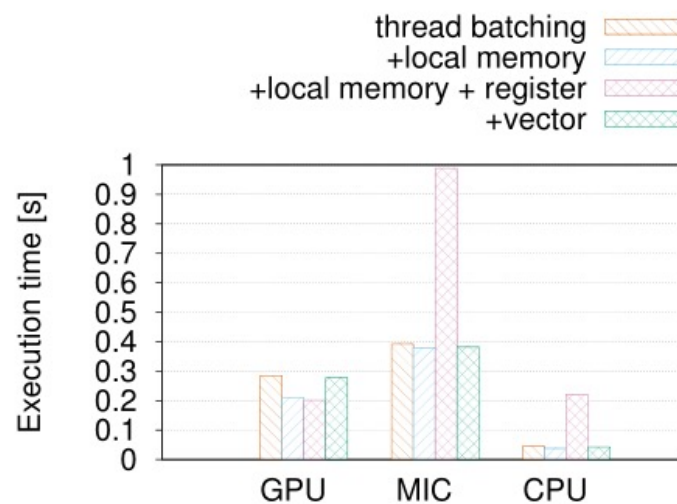
(a) Movielens



(b) Netflix



(c) YahooMusic R1



(d) YahooMusic R4

Goal:

Select the most appropriate

implementations for a

specific execution

context (i.e. target

architectures and

input datasets).

• 4.SETUP and DATASET •

Platform Configurations

- ✓ Intel Xeon E5-2670 (CPU): 16 cores
- ✓ NVIDIA Tesla K20c (GPU): 13 SM, 192 CUDA cores in each SM
- ✓ Intel Xeon Phi 31SP (MIC): 57 cores, 6GB global memory
- ✓ OpenCL (version 1.2)
- ✓ Host CPU: Redhat (v7.0) GCC(v4.9.2)

Input Datasets

■ Format of datasets

<user ID, item ID, rating>

	Abbr.	m	n	Training N_z
Movielens10M	MVLE	71567	65133	8000044
Netflix	NTFX	480189	17770	99072112
YahooMusic R1	YMR1	1948882	98212	115248575
YahooMusic R4	YMR4	7642	11916	211231

5. PERFORMANCE RESULTS



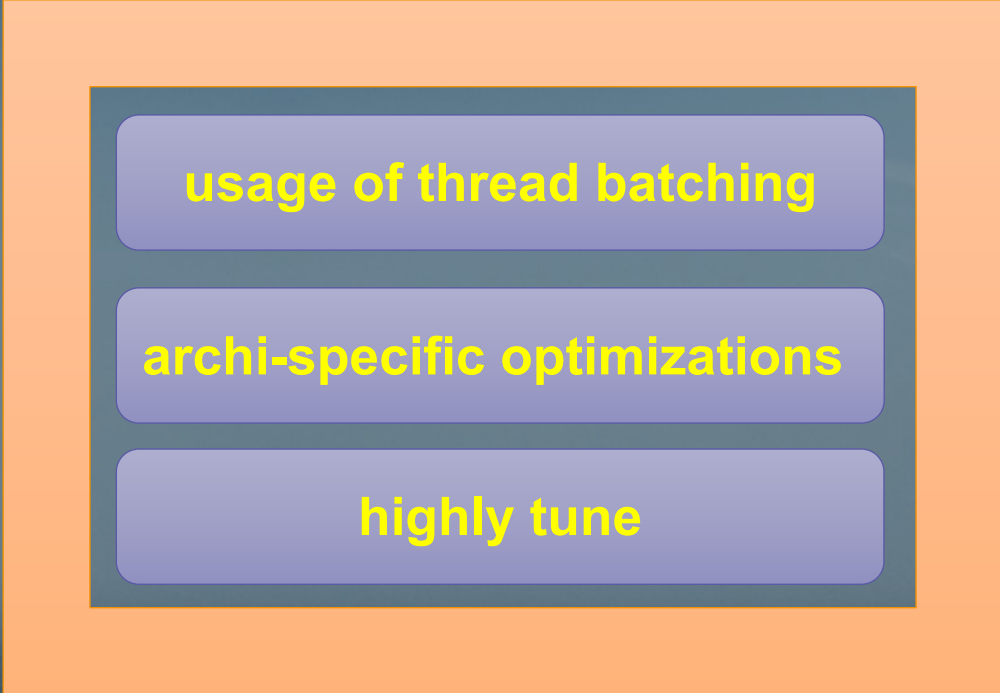
5.1 Compare with the State-of-the-art

■ vs **SAC15** { 5.5× faster than OpenMP
21.2× faster on K20c GPU

■ vs **HPDC16** (CuMF) 2.2× ~ 6.8×



Using cusparse library



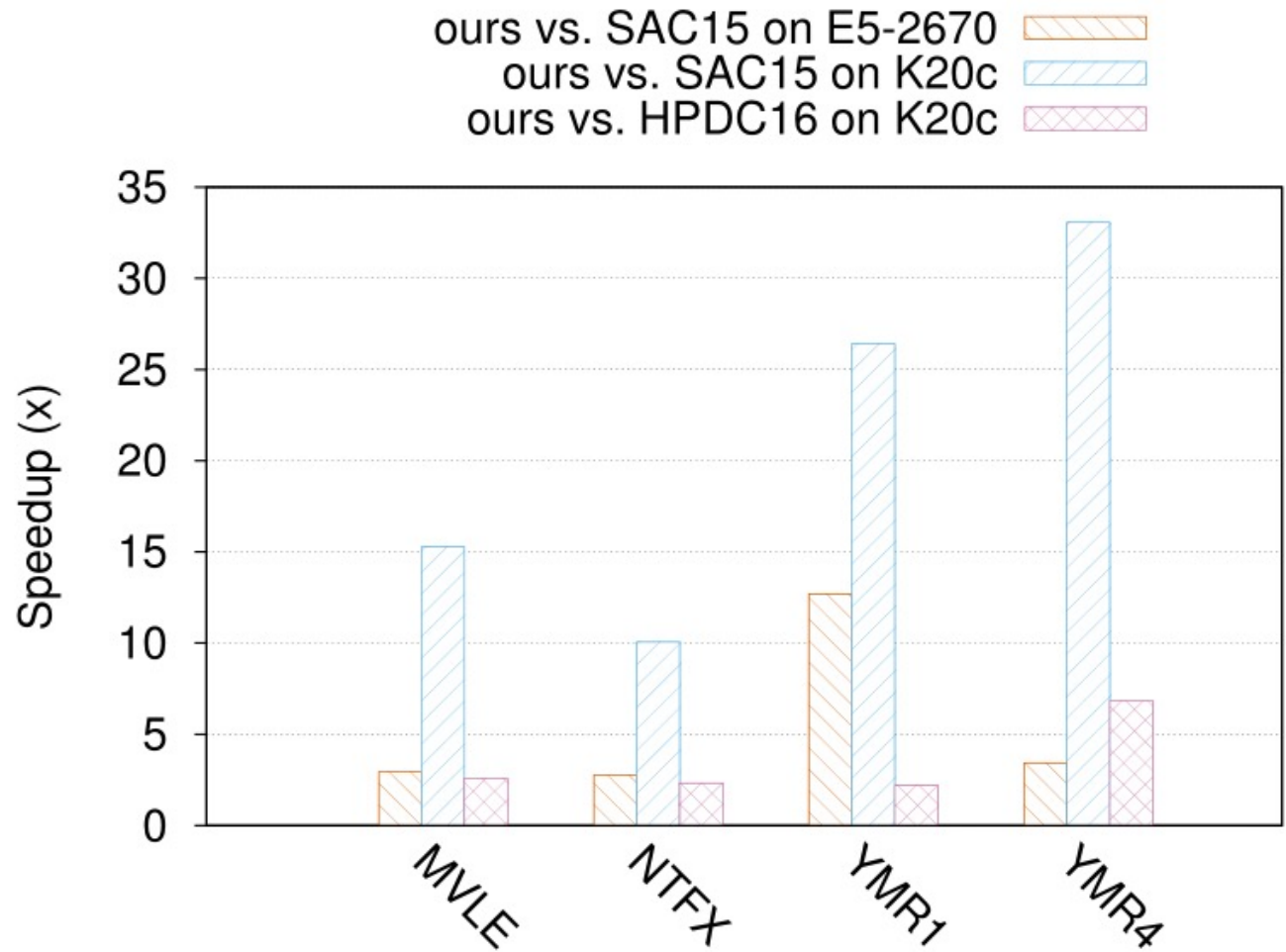
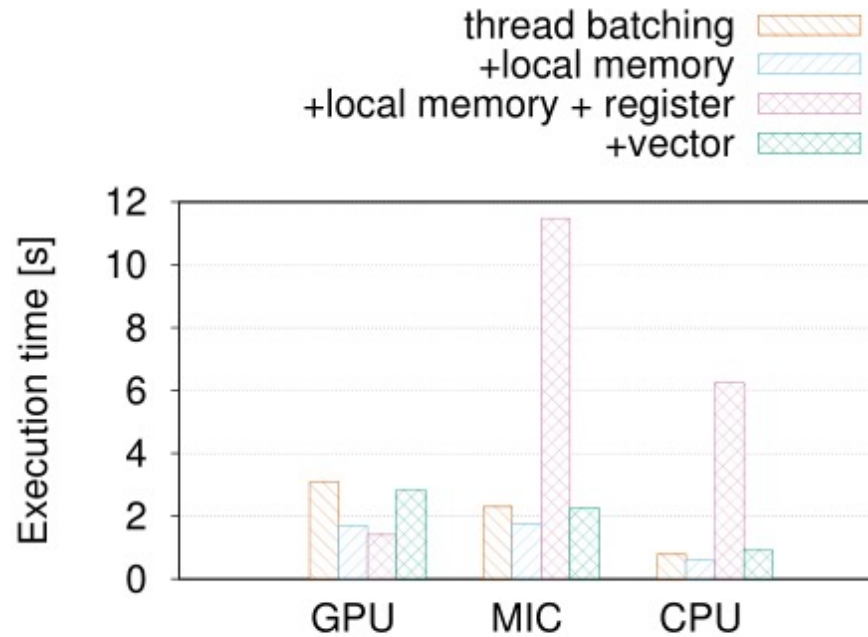
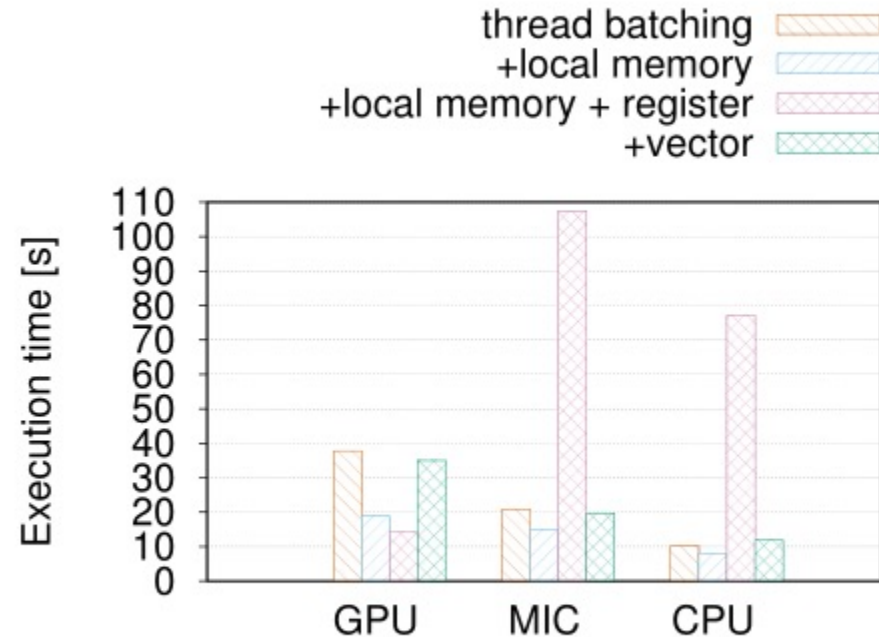


Figure 3. A performance comparison of our implementation versus the state-of-the-art implementations.

5.2 Evaluate Optimizations



(a) Movielens



(b) Netflix

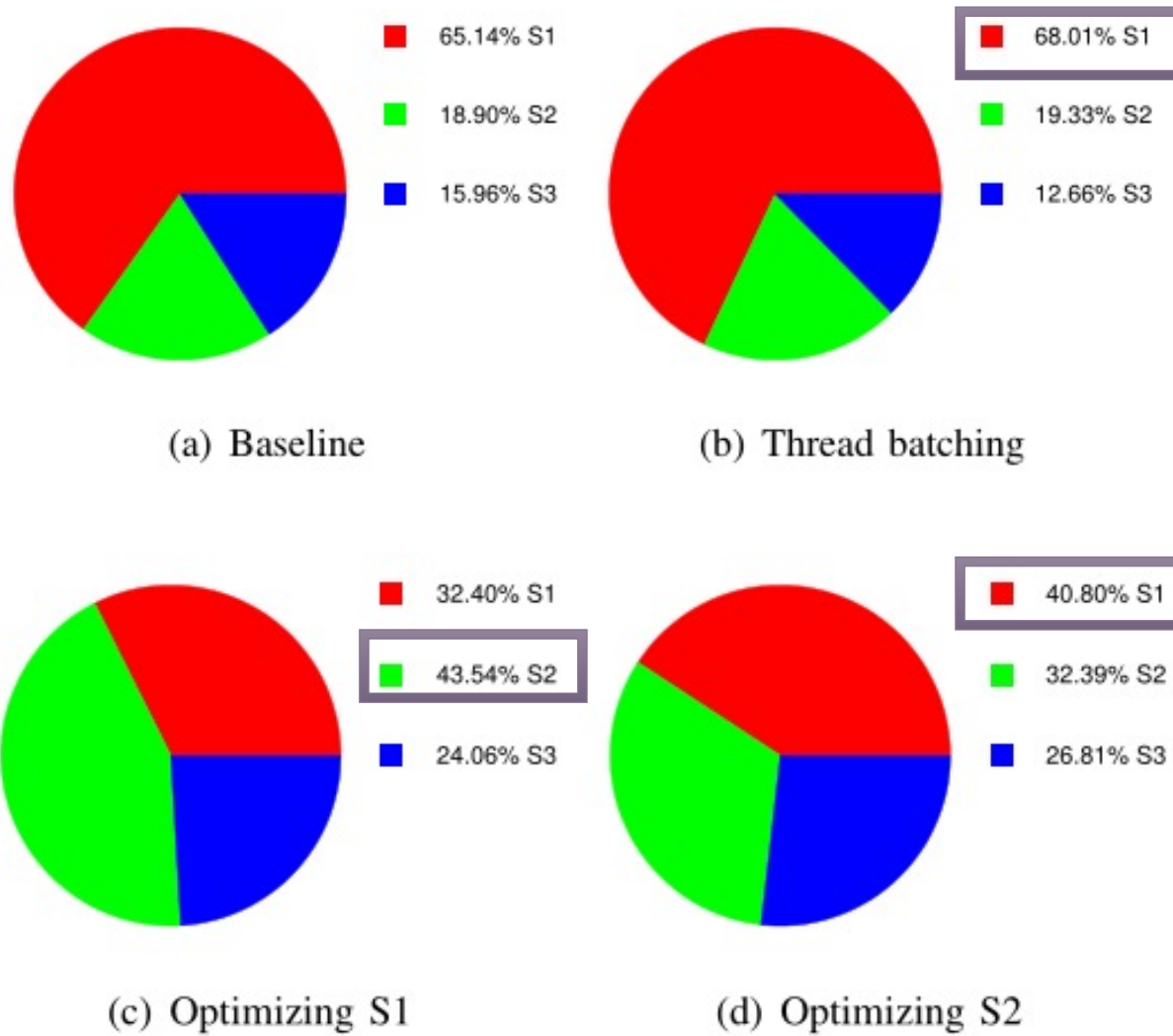
GPU: thread batching + local memory + register (upto $2.6\times$)

CPU/MIC: thread batching + local memory ($1.4\times$ for MIC, $1.6\times$ for CPU)

5.3 Apply Optimizations

- ◆ ALS {
 - step1: $Y^T Y + \lambda I$
 - step2: $Y^T r_u$
 - step3: solve the linear system
- ◆ give a priority to the most time-consuming step

Applying Optimizations



step a



+ thread batching

step b



optimize $Y^T Y + \lambda I$

step c

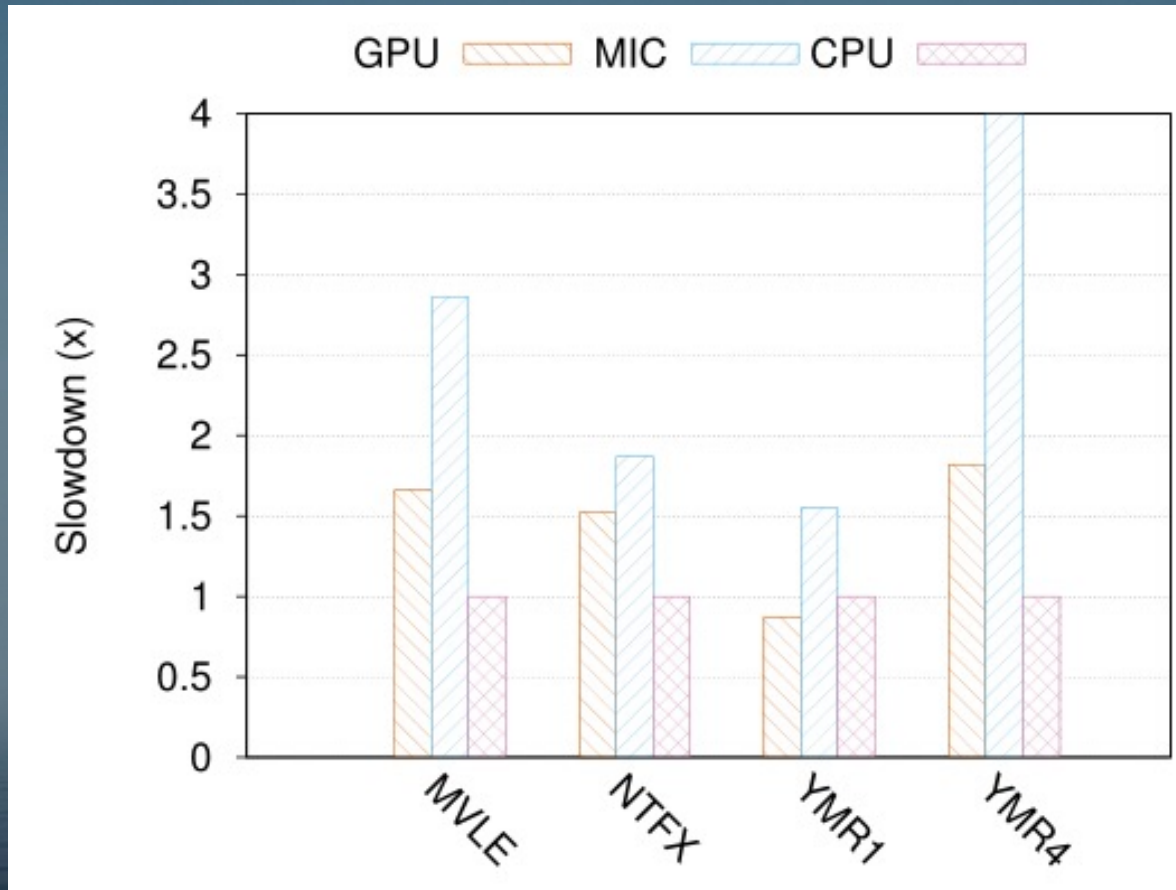


optimize $Y^T r_u$

step d

Hot-spot guided manner

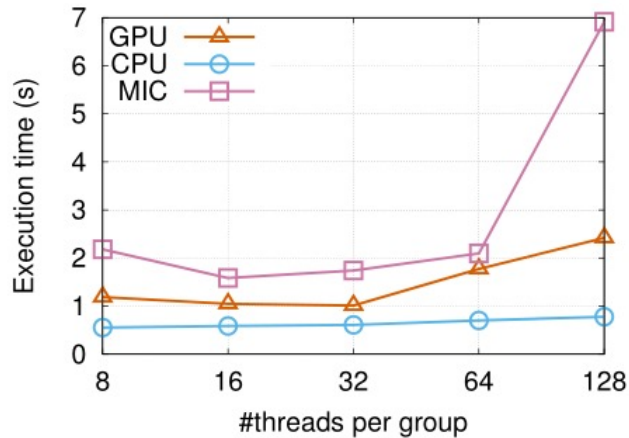
5.4 Compare between Different Architectures



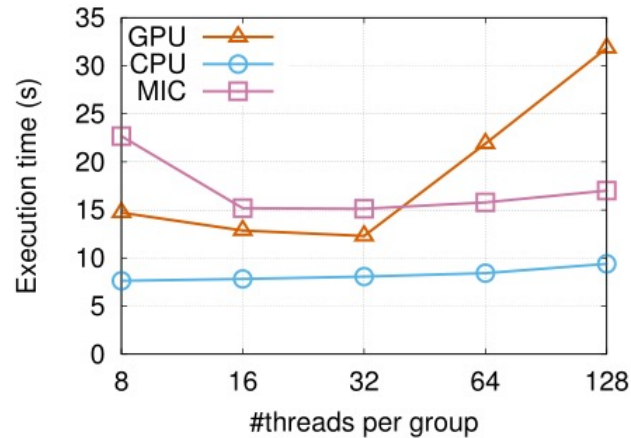
- CPU performs best
- 1.5x slower on GPU
- 4.1x slower on MIC

But, for Yahoo Music R1 the performance on GPU outperform that on 16-core CPU

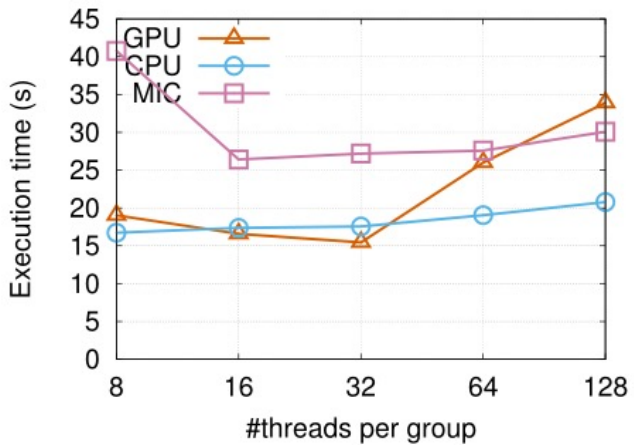
5.5 Sensitivity to Thread Blocks



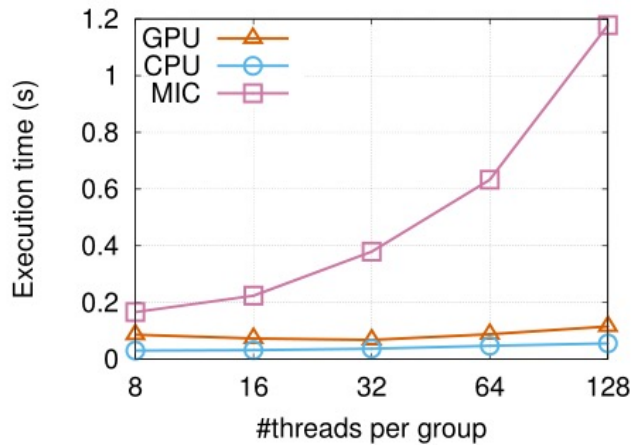
(a) Movielens



(b) Netflix



(c) YahooMusic R1

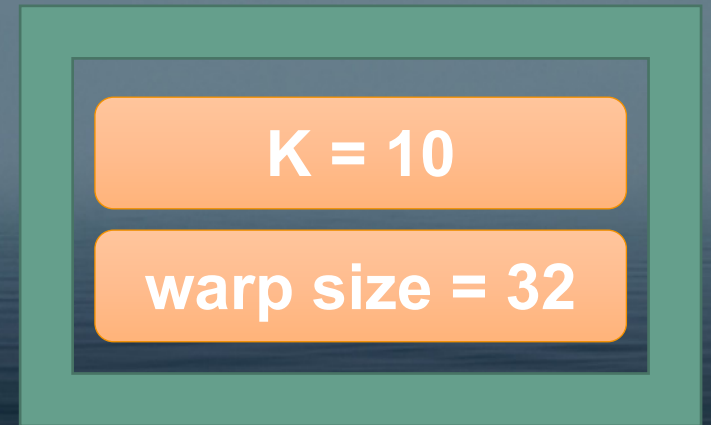


(d) YahooMusic R4

Configuration:

- 8192 * 32, k=10
- thread batching + local memory + registers

GPU: threads per block=16 / 32, best performance!



• 6. CONCLUSION •

Efficient and Portable ALS solver

- ✓ hierarchical thread organization on modern hardware
- ✓ thread batching
- ✓ architecture-specific optimizations
- ✓ OpenCL implementation (CPUs, GPUs, MICs)
- ✓ select suitable variant for each platform

THANK YOU