

MASCOT: A Quantization Framework for Efficient Matrix Factorization in Recommender Systems



Yunyong Ko¹, Jae-Seo Yu¹, Hong-Kyun Bae¹, Yongjun Park¹, Dongwon Lee², and Sang-Wook Kim¹

Hanyang University, Republic of Korea¹

The Pennsylvania State University, PA, USA²



□ Background

- Quantization & precision switching
- Matrix factorization (MF)

□ Proposed framework: MASCOT

- Motivation
- Strategy 1: m -quantization
- Strategy 2: g -switching

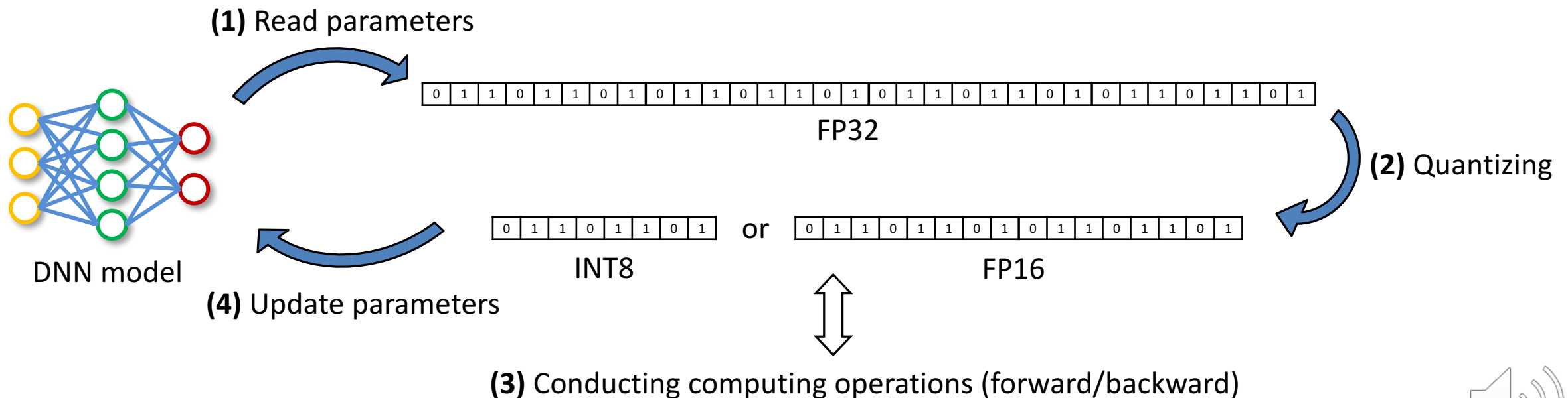
□ Experiments

□ Conclusions



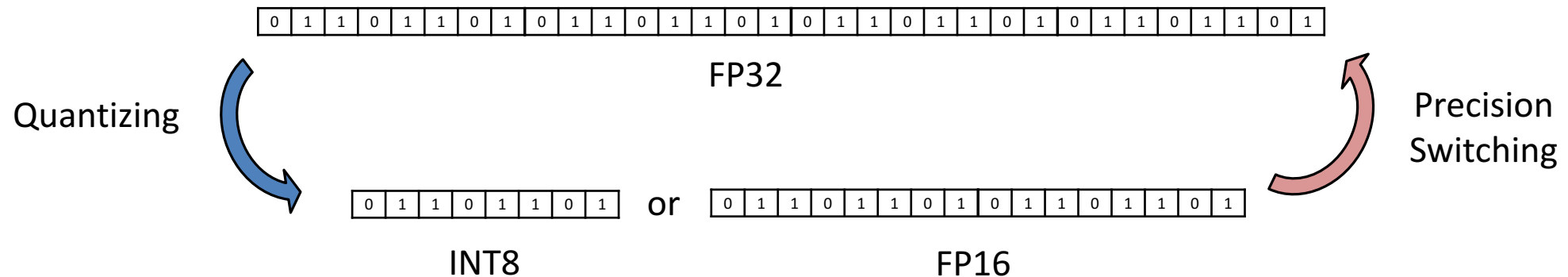
□ Quantization

- Converting the parameter value into a lower precision
 - Ex) FP32 (i.e., single precision) → FP16 (i.e., half precision)
- For improving the training performance of DNN models
 - To reduce *the overhead of computing operations* and *the memory usage*



□ Precision switching

- The quantization error may result in **degrading the model quality**
 - Quantization error: the error caused by low precision
- Switching back the low precision to high precision to prevent the loss of accuracy
 - i.e., FP16 (i.e., half precision) → FP32 (i.e., single precision)



□ Matrix Factorization (MF)

- One of the popular *collaborative filtering* algorithms in recommender systems (RS)
- Aiming to obtain the latent matrices \mathbf{P} and \mathbf{Q} (satisfying $\mathbf{R} \approx \mathbf{P} \times \mathbf{Q}^T$)

$$\begin{matrix} & \overbrace{\hspace{2cm}}^n \\ \underbrace{\hspace{1cm}}_m \left\{ \begin{array}{ccccc} 5 & 2 & \text{ } & \dots & 1 \\ 4 & \text{ } & 4 & \dots & \text{ } \\ 3 & 3 & \text{ } & \dots & \text{ } \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 5 & \text{ } & \text{ } & \dots & 3 \end{array} \right\} & \approx & \begin{matrix} \overbrace{\hspace{2cm}}^k \\ \begin{array}{ccc} 2.7 & \dots & 0.5 \\ 0.9 & \dots & 0.3 \\ 0.9 & \dots & 1.5 \\ \vdots & \ddots & \vdots \\ 1.9 & \dots & 0.2 \end{array} \end{matrix} & \times & \begin{matrix} \overbrace{\hspace{2cm}}^n \\ \begin{array}{ccccc} 0.3 & 0.7 & 1.5 & \dots & 1.2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1.6 & 0.5 & 0.1 & \dots & 2.5 \end{array} \end{matrix} & = & \begin{matrix} & \overbrace{\hspace{2cm}}^n \\ \underbrace{\hspace{1cm}}_m \left\{ \begin{array}{ccccc} 4.9 & 1.5 & 3.5 & \dots & 1 \\ 3.5 & 2.4 & 4.0 & \dots & 1.9 \\ 3.0 & 3.5 & 4.5 & \dots & 0.5 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 4.8 & 1.3 & 3.4 & \dots & 3.2 \end{array} \right\} \end{matrix} \\ & \mathbf{R} & & \mathbf{P} & & \mathbf{Q}^T & & \hat{\mathbf{R}} \end{matrix}$$

□ Challenge

- The growing scale of users/items and model architectures can significantly **slow down the training of MF models**

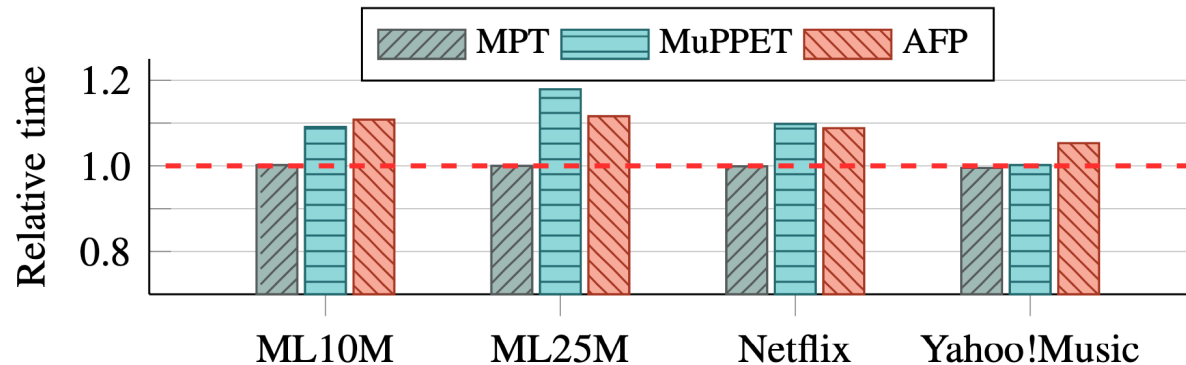


□ Our question

- “Can the proven quantization be adopted to improve the training of MF models in recommender systems?”

□ Preliminary experiment with four RS datasets

- SOTA quantization methods are *rarely effective* in MF model training



MPT: P. Micikevicius et al. Mixed precision training. In *ICLR*, 2018.

MuPPET: A. Rajagopal et al., MuPPET: A precision-switching strategy for quantised fixed-point training of CNNs. In *ICML*, 2020.

AFP: X. Zhang et al. Fixed-point back-propagation training. In *CVPR*, 2020.



Analysis on the Training of a MF Model

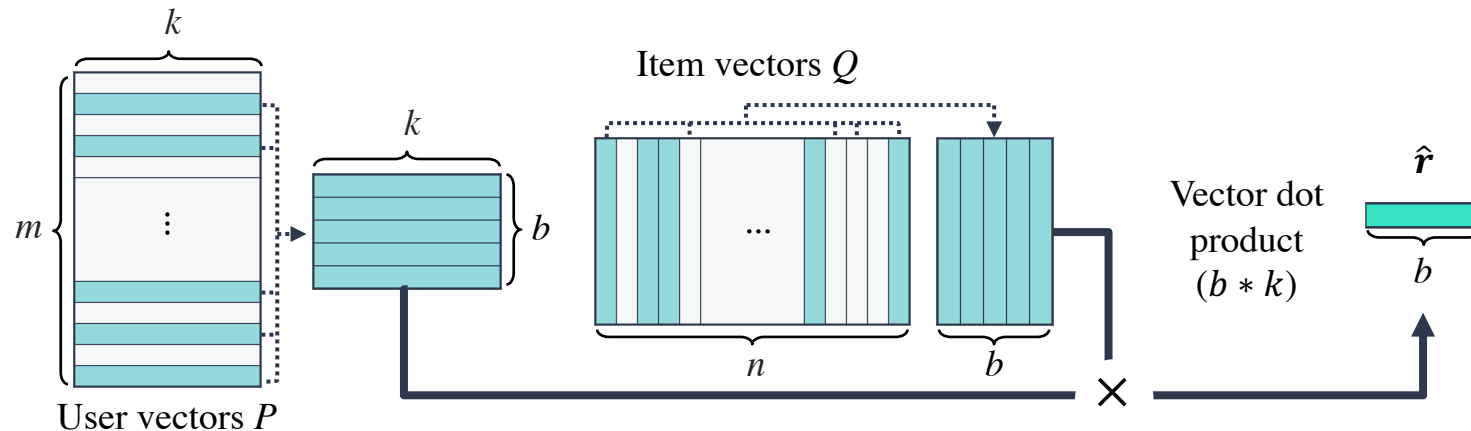


□ Training of a MF model is much more *memory-intensive* than that of a DNN model (*Observation 1*)

■ Memory and computation costs of a MF model (vector dot product)

□ Memory cost ($b * 2k$) $>$ Computation cost ($b * k$)

■ *Little room* for the performance improvement by the quantization



<Computational cost of a MF model>

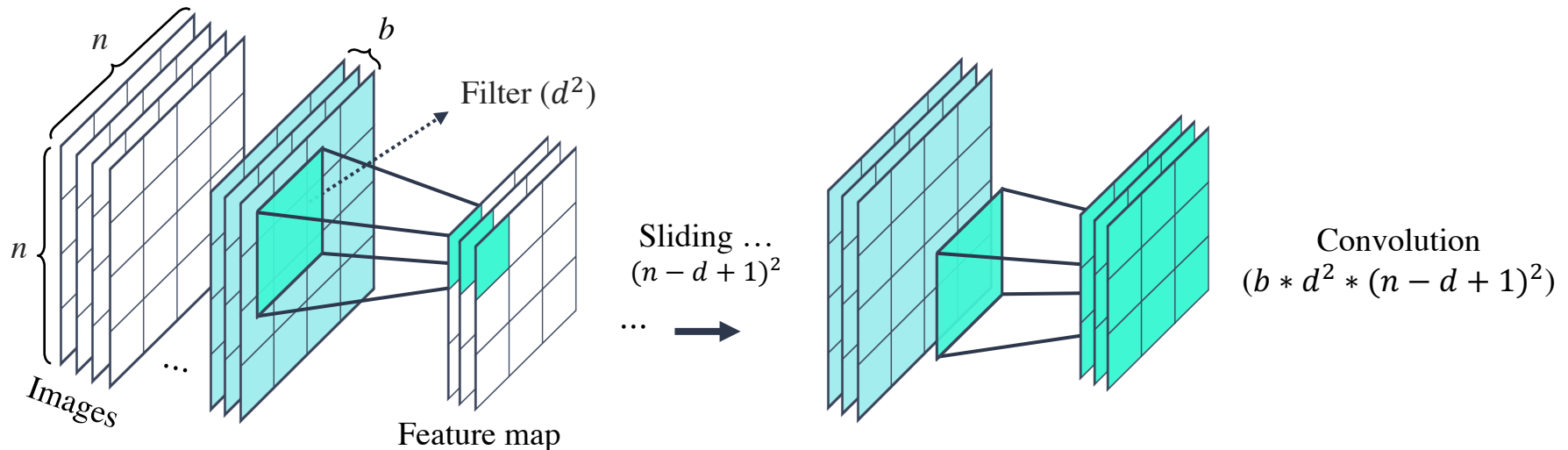


□ The DNN model training is *computation-intensive*

■ Memory and computation costs of a convolution layer

□ Computation cost ($b * d^2 * (n - d + 1)^2$) > Memory cost ($b * n^2 + d^2$)

■ The room for the performance improvement by the quantization is sufficient



<Computational cost of a convolution layer>



Unique Feature of Datasets in RS

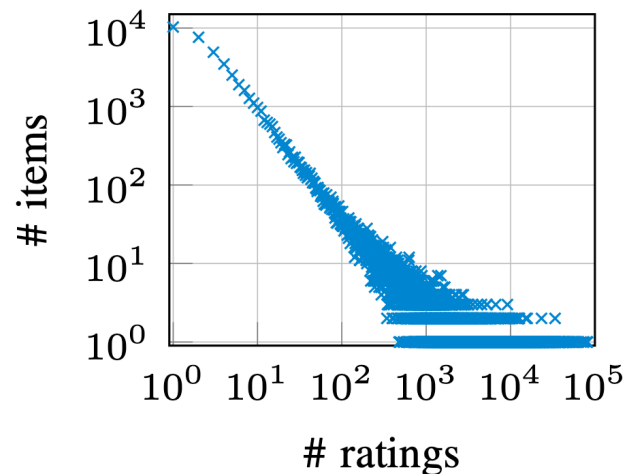


□ The datasets used in RS have a *power-law distribution*

- A *majority* of users/items have a *small* number of ratings
- A *small* number of users/items have a very *large* number of ratings

□ The MF model training with SGD

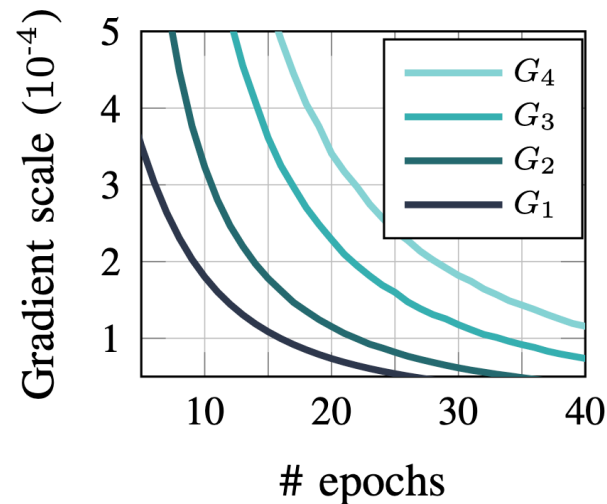
- The update frequency for each latent vector depends on the number of ratings
 - The latent vectors of *a few* users/items with many ratings are updated *frequently*
 - The latent vectors of *many* users/items with few ratings are updated *infrequently*



Quantization Error of Each Latent Vector



- The scale of the gradient for each latent vector *varies*, depending on the number of ratings
 - The scale of gradient tends to decrease as the model is trained more
- The quantization error of each user/item is *different* from each other, depending on the number of ratings (*Observation 2*)



- Considering the quantization error for the entire model (i.e., not considering the difference among users/items)
 - The loss of accuracy
 - For users with many ratings, the precision switching is likely to be applied *too late*
 - The reduced training performance
 - For users with few ratings, the precision switching is likely to be applied *unnecessarily quickly*



Overview of The Proposed Framework (MASCOT)



- Quantization strategy for memory access (***m*-quantization**)
 - Storing and managing the parameters of MF models in low precision

- Group-based precision switching strategy (***g*-switching**)
 - Grouping users/items having a similar number of ratings
 - Applying precision switching in a group-wise manner

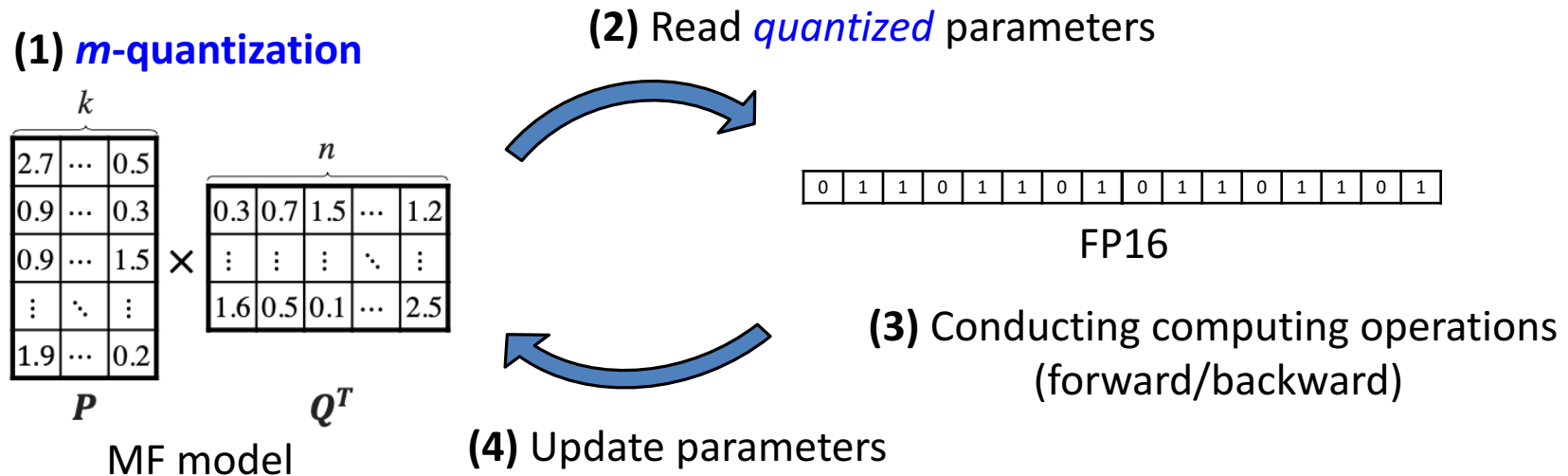


Strategy 1: m -quantization



□ Storing and managing the parameters of a MF model in low precision

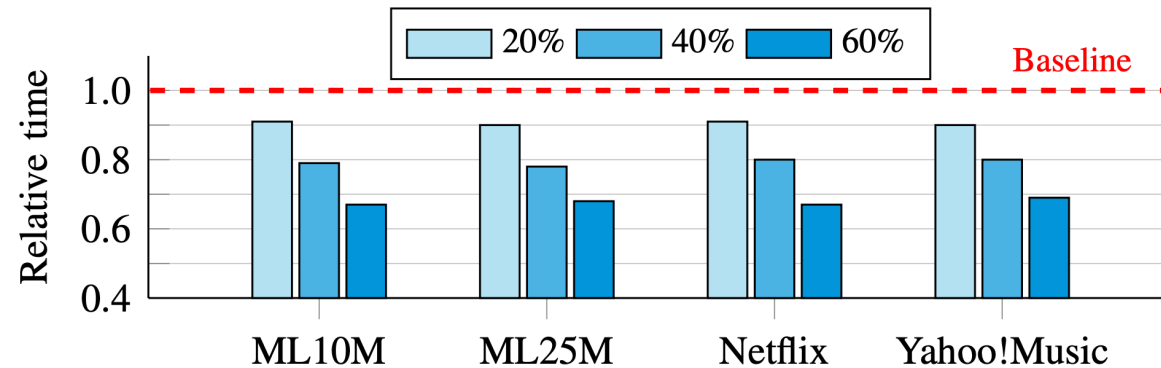
- For *improving the memory access operations* in the training of MF models



Strategy 1: m -quantization



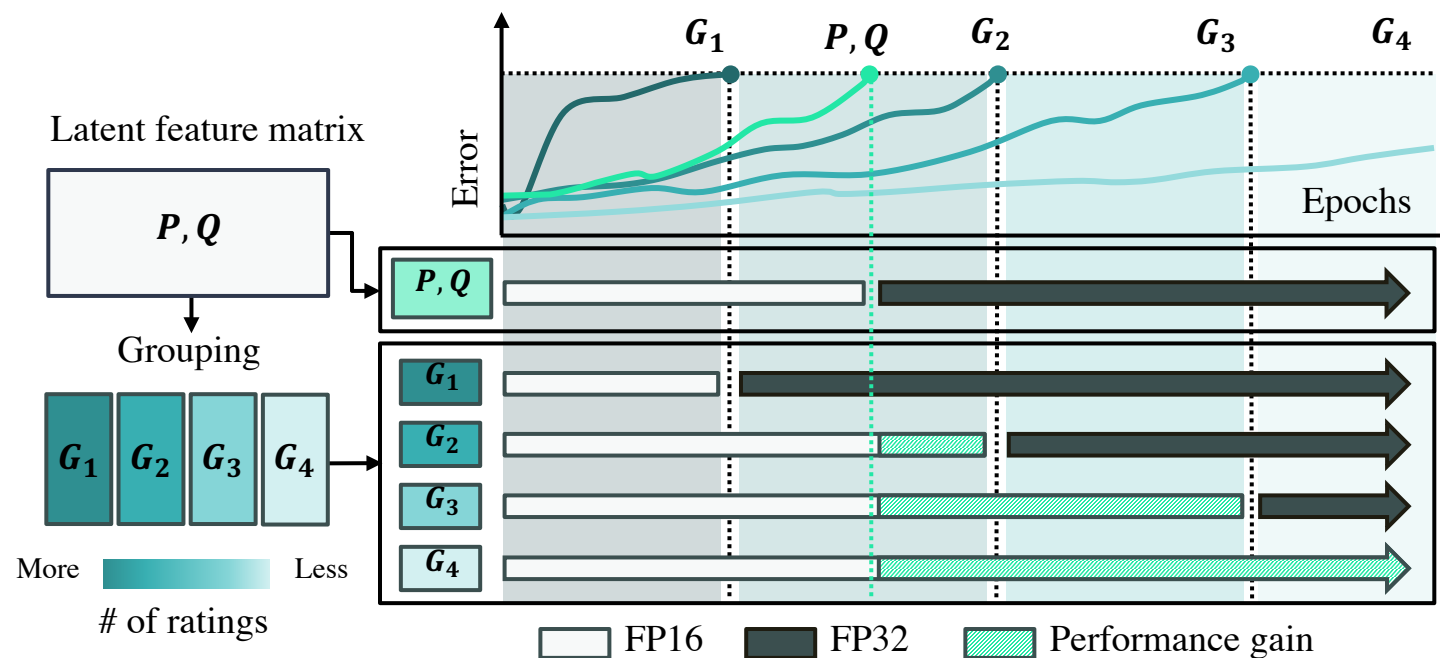
- The potential of the m -quantization strategy for improving the training performance of MF models



Strategy 2: *g*-switching

□ *Grouping* users/items having a *similar* number of ratings and *applying* precision switching in a *group-wise* manner

■ For *considering the difference among users/items* in RS datasets



<The performance improvement by the *g*-switching of MASCOT>

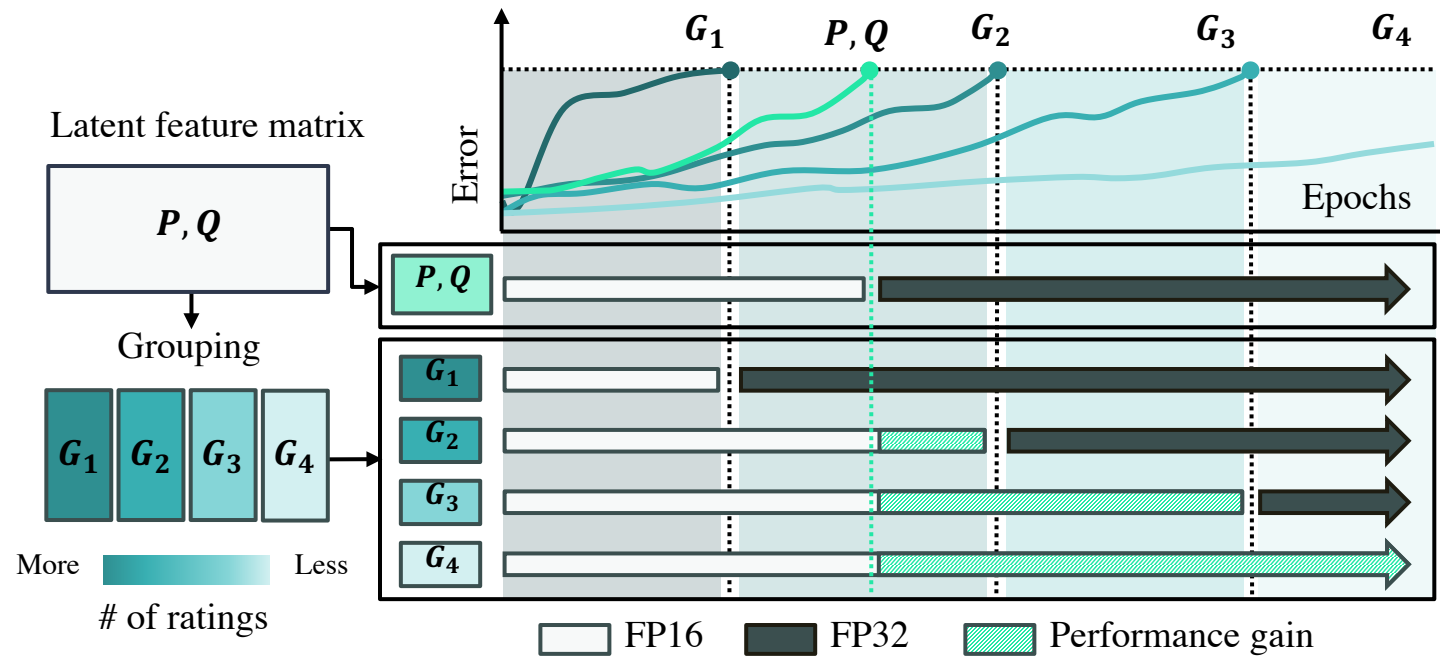


Strategy 2: *g*-switching



□ Performance improvement by the *g*-switching

- Applying precision switching *only to the groups highly likely to cause the model error*



<The performance improvement by the *g*-switching of MASCOT>



□ Training process of MASCOT

- 1) Grouping users/items vectors
- 2) For each user/item,
 - Read the embedding vector;
 - Compute the loss and gradients;
 - Update the embedding vector;
 - Compute the quantization error of each group;
- 3) Determining to apply precision switching to each user/item group

Algorithm 1 Training of MASCOT

Require: $R \in \mathbb{R}^{m \times n}$, $P \in \mathbb{R}^{m \times k}$, $Q \in \mathbb{R}^{n \times k}$, # of groups g , error estimate period π , sample ratio γ , error threshold θ , learning rate η

```
1:  $R, P, Q \leftarrow \text{grouping}(R, P, Q, g)$ 
2: Initialize  $P, Q$  with half precision
3: for  $t = 1, \dots, T$  do
4:   for each rating  $r_{u,i}$  do
5:      $\nabla p_u \leftarrow e_{u,i} \cdot q_i - \lambda_P \cdot p_u$ 
6:      $\nabla q_i \leftarrow e_{u,i} \cdot p_u - \lambda_Q \cdot q_i$ 
7:      $p_u \leftarrow p_u + \eta \cdot \nabla p_u$ 
8:      $q_i \leftarrow q_i + \eta \cdot \nabla q_i$ 
9:     if  $S \sim B(\gamma)$  then
10:       $S_j^U \cdot \text{push}(\nabla p_u), S_j^I \cdot \text{push}(\nabla q_i)$ 
11:    end if
12:     $\text{update\_latent\_matrix}(P, Q, p_u, q_i)$ 
13:  end for
14:  if  $t \pmod{\pi} == 0$  then
15:    for  $j = 1, \dots, g$  do
16:       $\epsilon^U \leftarrow q\text{-error}(S_j^U), \epsilon^I \leftarrow q\text{-error}(S_j^I)$ 
17:       $\text{precision\_switching}(P, Q, \epsilon^U, \epsilon^I, \theta)$ 
18:    end for
19:  end if
20:   $\forall j \in \{1, \dots, g\}, S_j^U \cdot \text{flush}(\cdot), S_j^I \cdot \text{flush}(\cdot)$ 
21: end for
22: Return  $P, Q$ 
```



□ Models & Datasets

- MF models with varying the dimensionality of laten space (64, 128)
- Statistics of datasets

Datasets	# of users	# of items	# of ratings	Sparsity
ML10M	69,878	10,677	10,000,035	98.66%
ML25M	162,541	59,047	24,997,208	99.74%
Netflix	480,189	17,770	100,480,507	98.82%
Yahoo!Musics	1,000,990	624,961	256,804,235	99.96%

□ Competing algorithms

- MPT (ICML'18)
- MuPPET (ICML'20)
- AFP (CVPR'20)
- Two baselines (FP32, FP16)



□ Q1: Training performance

- Does MASCOT improve the training performance of MF models more than existing quantization methods?

□ Q2: Model quality

- Does MASCOT provide the errors of MF models lower than existing quantization methods?

□ Q3: Effectiveness of each strategy

- How effective are the strategies of MASCOT in improving the MF model training?



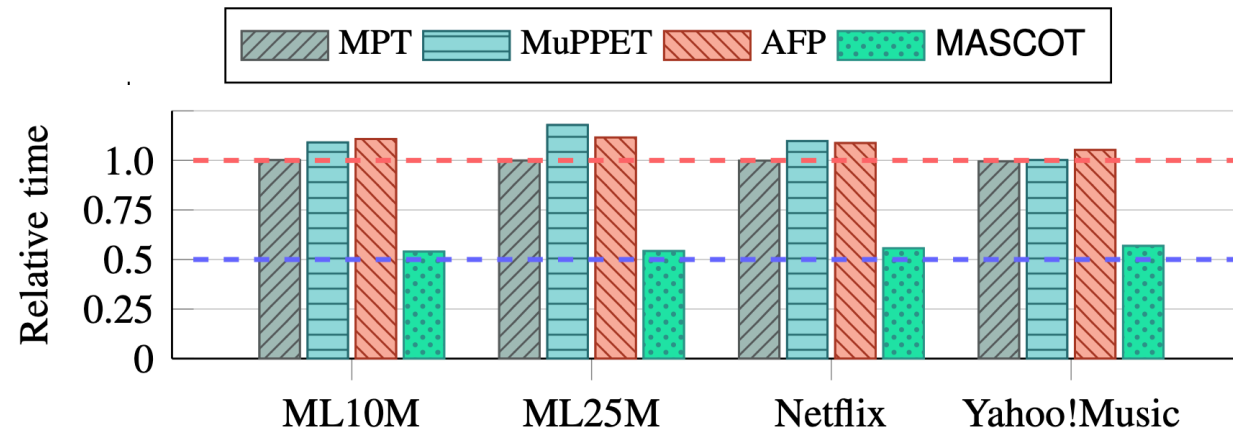
Q1. Training Performance



□ MASCOT improves the training performance of MF models most

■ About *45% performance improvement* on average

□ Existing SOTA quantization methods *degrade* the training performance of MF models

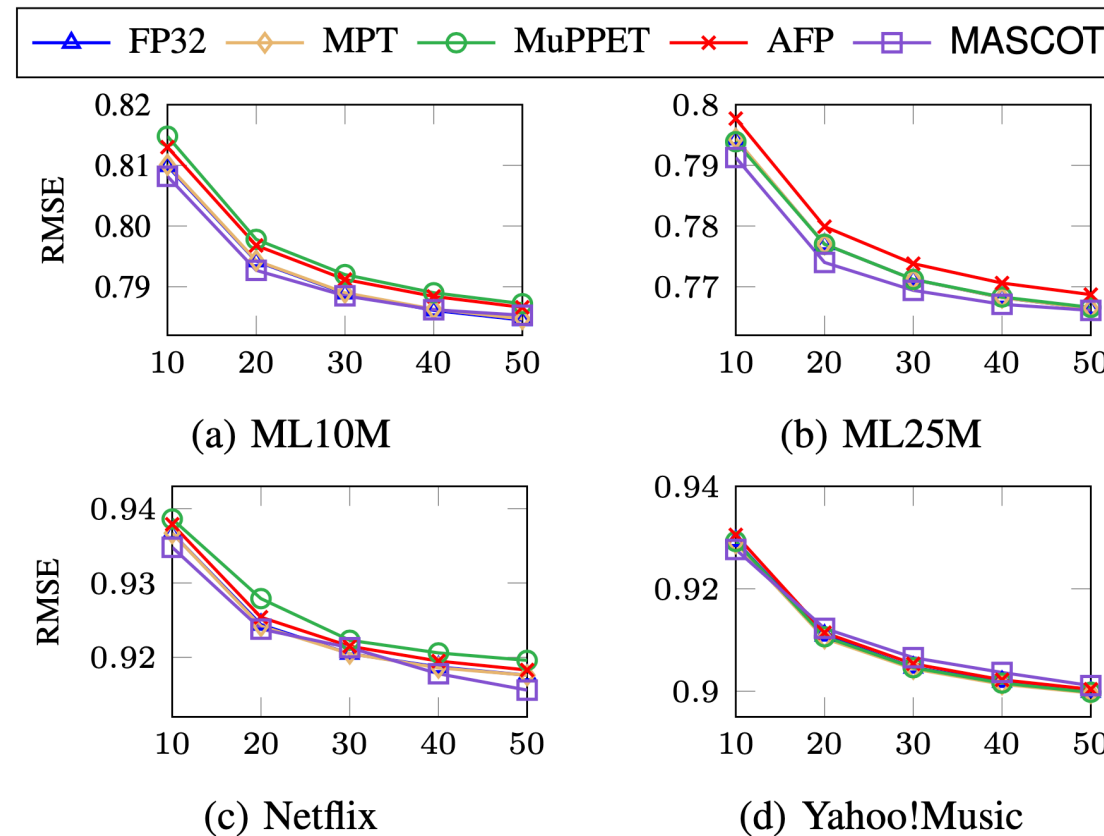


Q2. Model Quality



□ MASCOT achieves low model errors comparable to that of FP32

- The g-switching applies the precision switching only to the groups that are highly likely to incur significant model errors

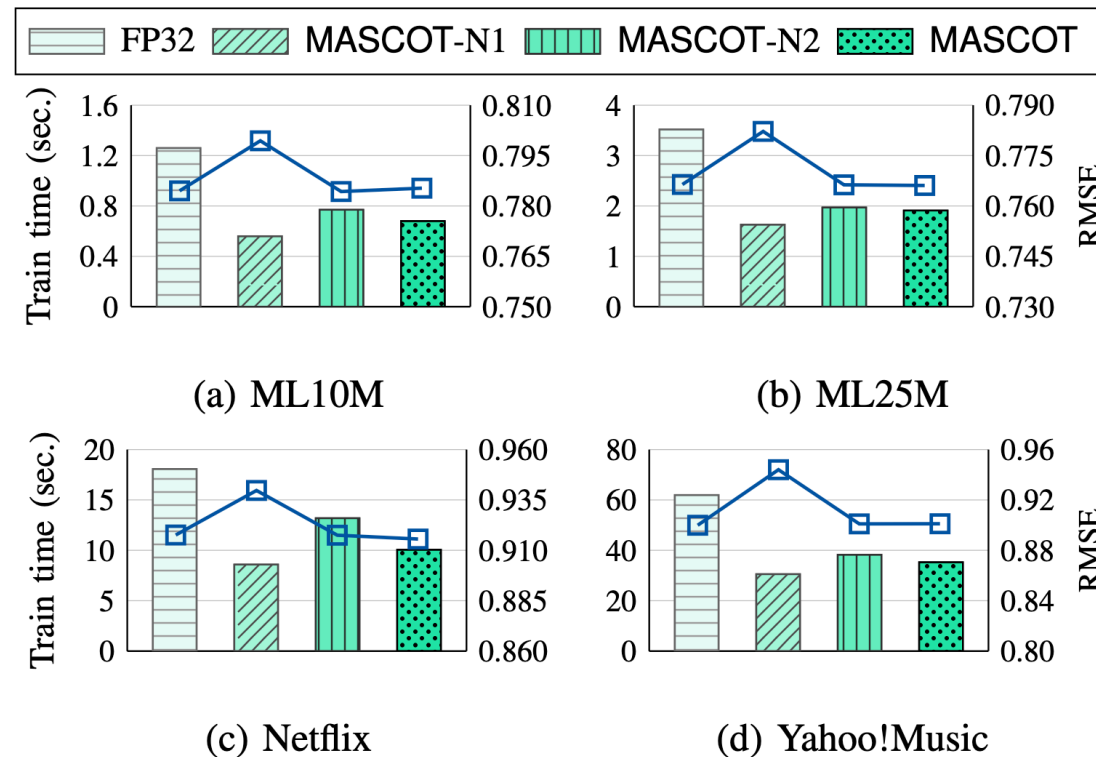


Q3. Ablation Study: Effectiveness of Strategies



Three versions of MASCOT

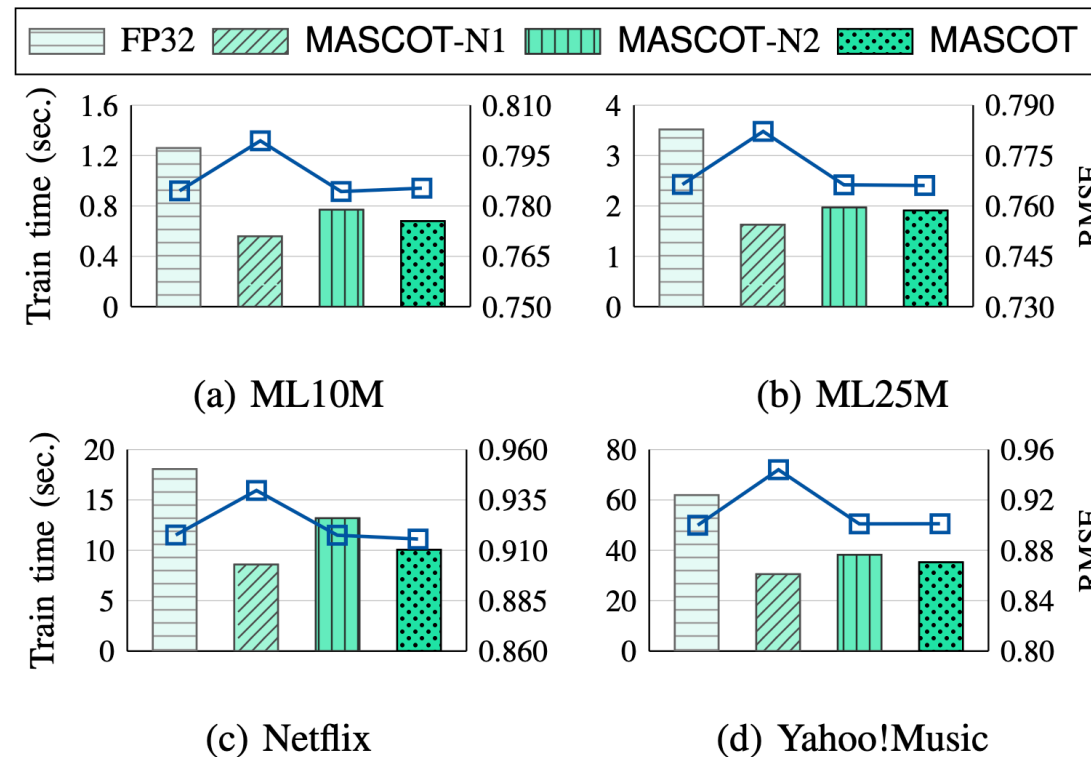
- MASCOT-N1 is with only *m*-quantization
- MASCOT-N2 is with *m*-quantization and the existing precision switching method
- MASCOT is with both *m*-quantization and *g*-switching



Q3. Ablation Study: Effectiveness of Strategies



- Both strategies of MASCOT are quite effective in the MF model training.
- MASCOT-N1 shows the best improvement in reducing training time (*m-quantization*)
- MASCOT outperforms MASCOT-N2 in terms of both training performance and model quality (*g-switching*)



- ❑ **Discovering that existing SOTA quantization techniques are rarely effective in the training of MF models**
- ❑ **Identifying two unique features of the training of MF models**
 - (i) The training of MF models is more memory-intensive than that of DNN models
 - (ii) The quantization error of each user/item differs, depending on the number of
- ❑ **Proposing a quantization framework for efficient training of MF models**
 - Employing two strategies to address the unique features of the training MF models
- ❑ **Comprehensive evaluation verifying the effectiveness of MASCOT in the training MF models**
 - Improving the training performance by about 45% on average (almost ideal)



Thank You !

Email: koyunyong@hanyang.ac.kr

