



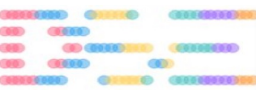
# Automated Machine Learning for Recommender Systems

Xiangyu Zhao

Data Science and Engineering Lab

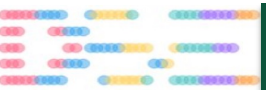
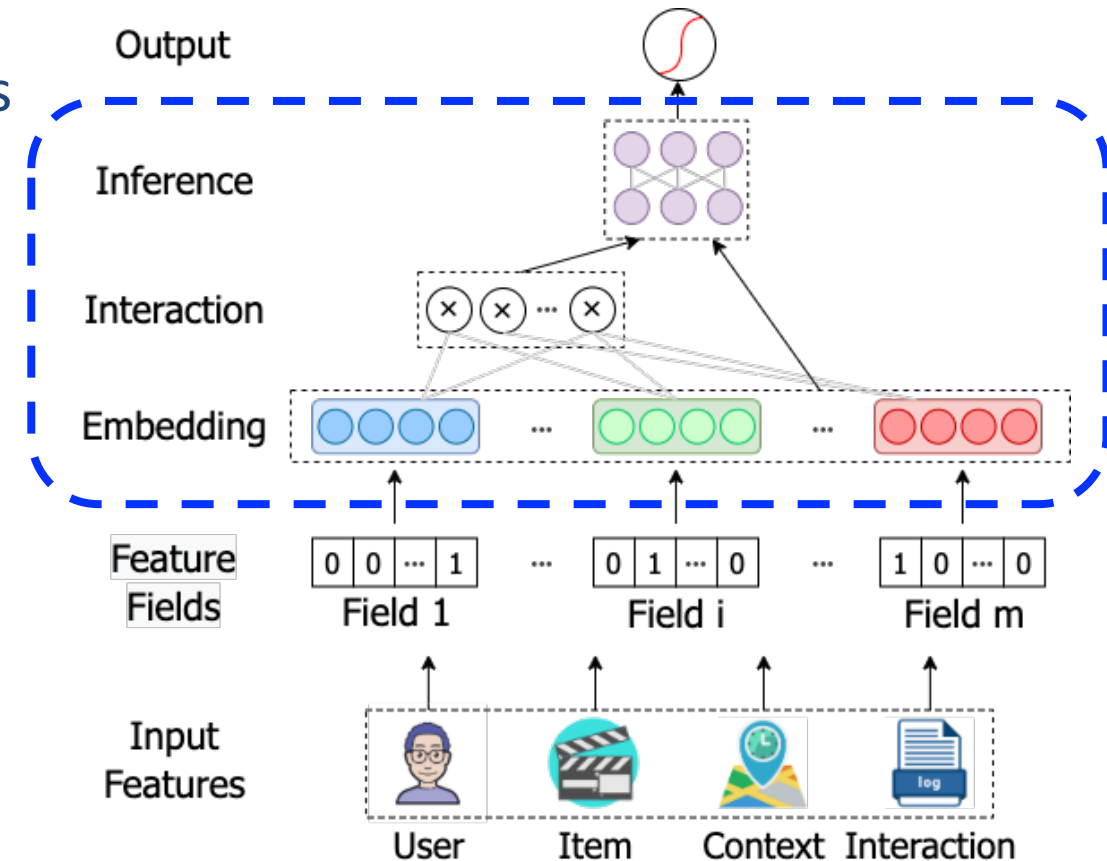
Michigan State University

[www.cse.msu.edu/~zhaoxi35](http://www.cse.msu.edu/~zhaoxi35) , [zhaoxi35@msu.edu](mailto:zhaoxi35@msu.edu)



# Deep Recommender Architectures

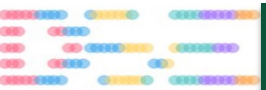
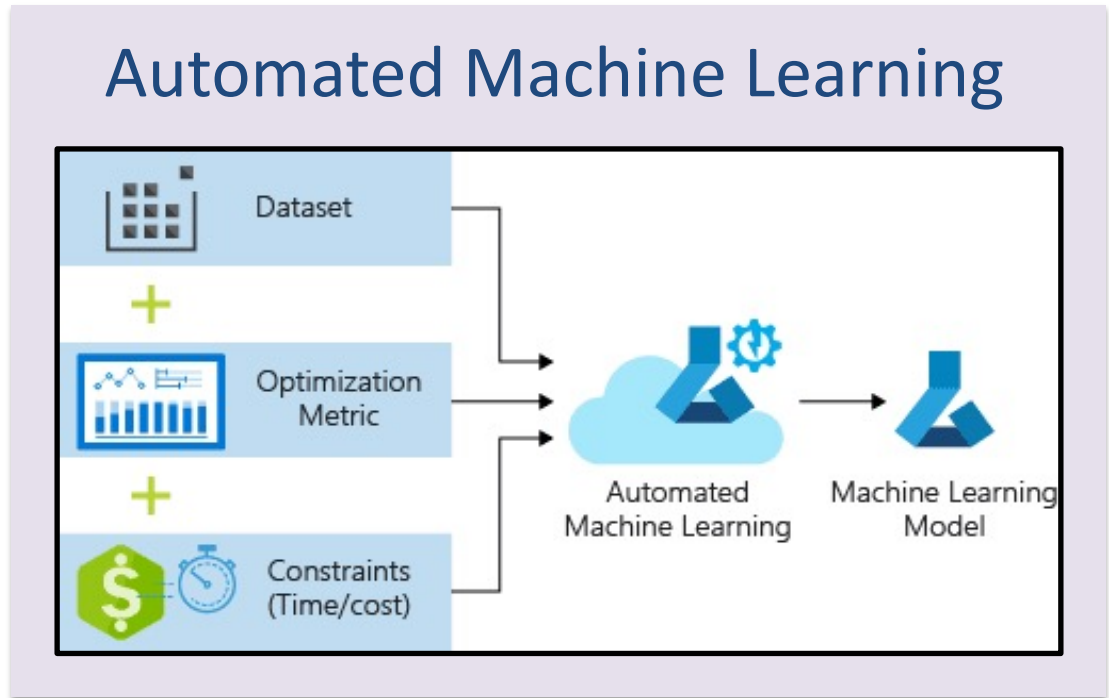
- Advantages
  - Feature representations of users and items
  - Non-linear relationships between users and items
- Typical architecture
  - Embedding layer
  - Interaction layer
  - Inference layer
- Manually designed architecture
  - Expert knowledge
  - Time and engineering efforts
  - Human error and bias → suboptimal architecture



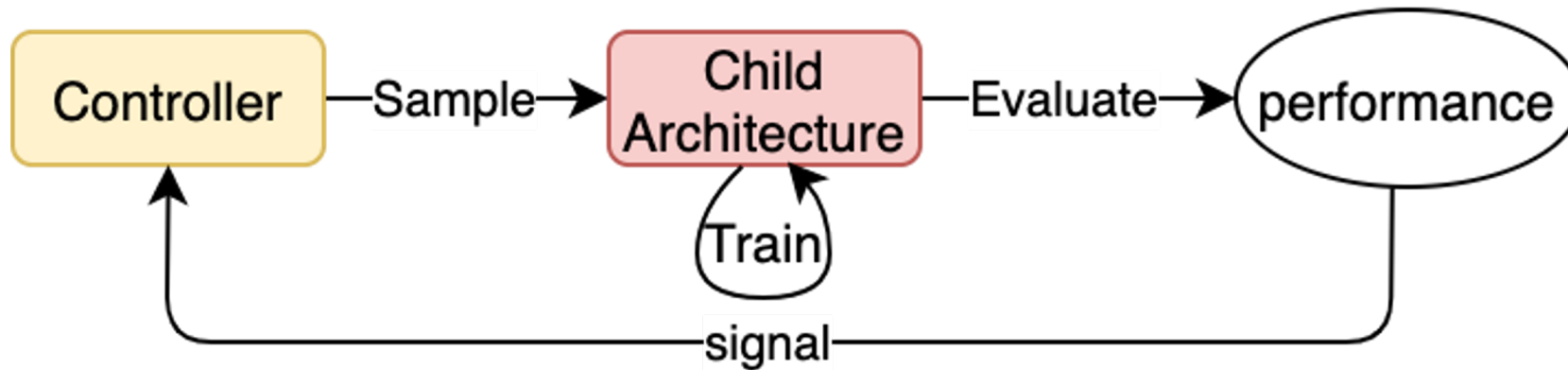
# AutoML for Deep Recommender Systems



- Deep architectures are designed by the machine automatically
- Advantages
  - Less expert knowledge
  - Saving time and efforts
  - Different data → different architectures



- Reinforcement Learning-based NAS
  - Controller: learning to select optimal child architecture
  - Child architecture: the DNN with a specific architecture



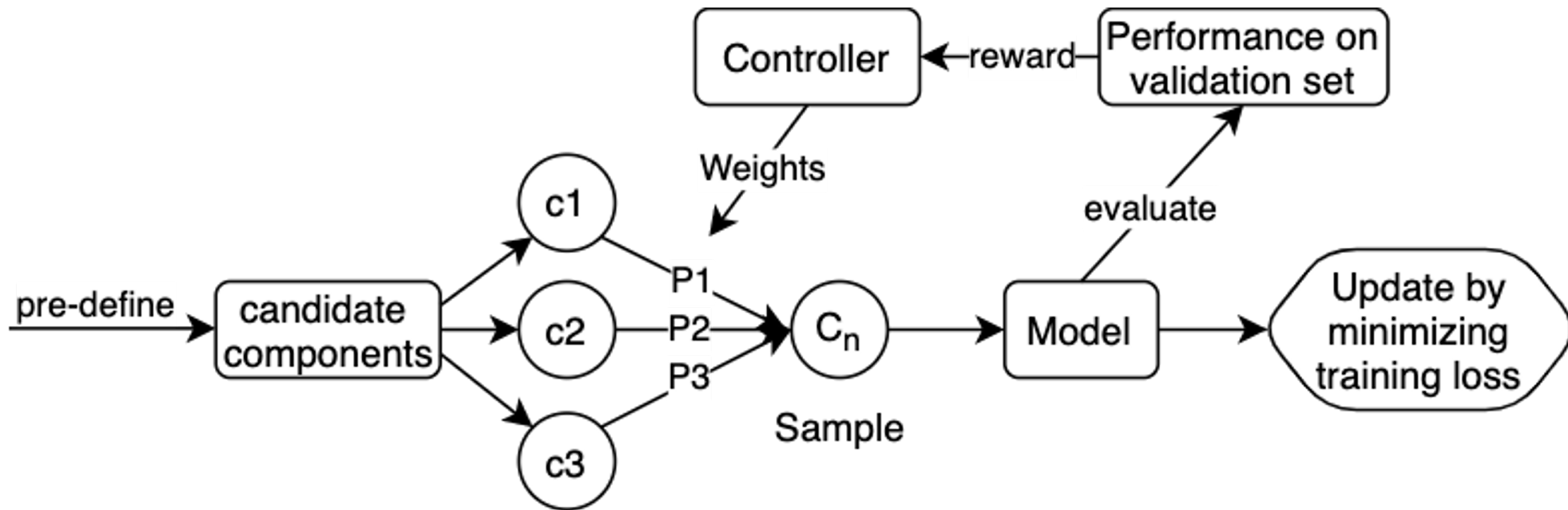
candidate  
components:

- identity
- 1x7 then 7x1 convolution
- 3x3 average pooling
- 5x5 max pooling
- 1x1 convolution
- 3x3 depthwise-separable conv
- 7x7 depthwise-separable conv
- 1x3 then 3x1 convolution
- 3x3 dilated convolution
- 3x3 max pooling
- 7x7 max pooling
- 3x3 convolution
- 5x5 depthwise-seperable conv



## Reinforcement Learning-based NAS

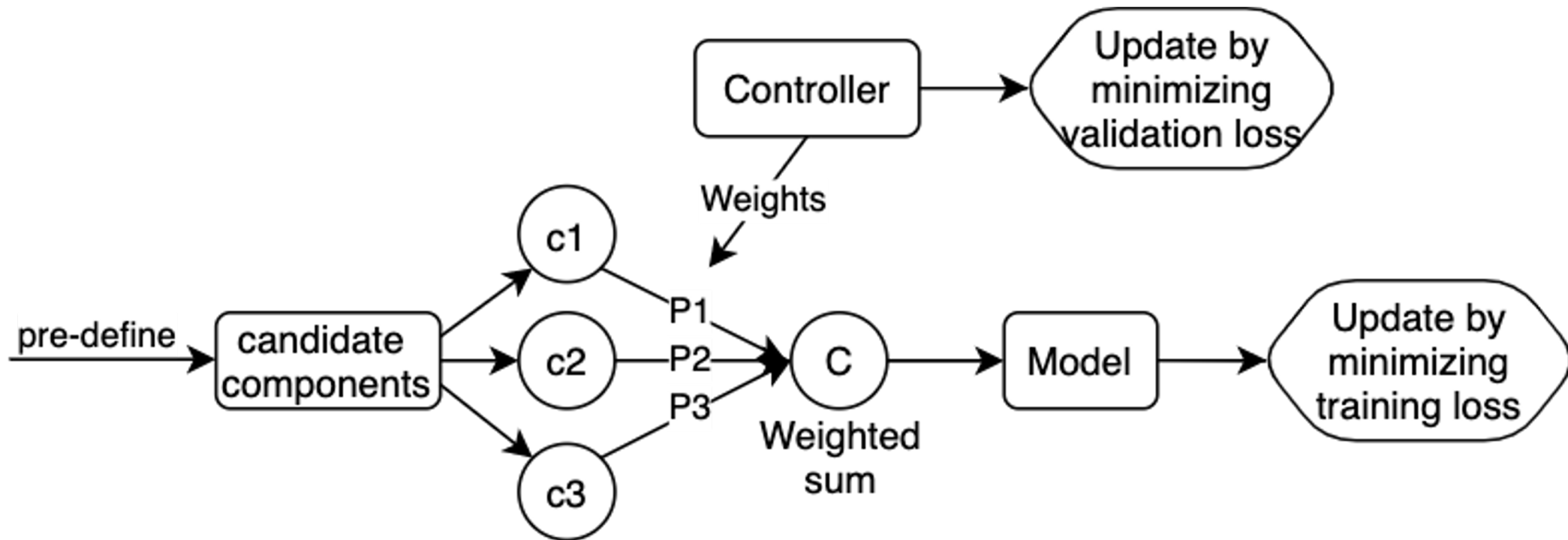
- Hard selection on candidate components
- The model's performance on validation set are viewed as reward
- The weights of controller are updated to maximize the reward



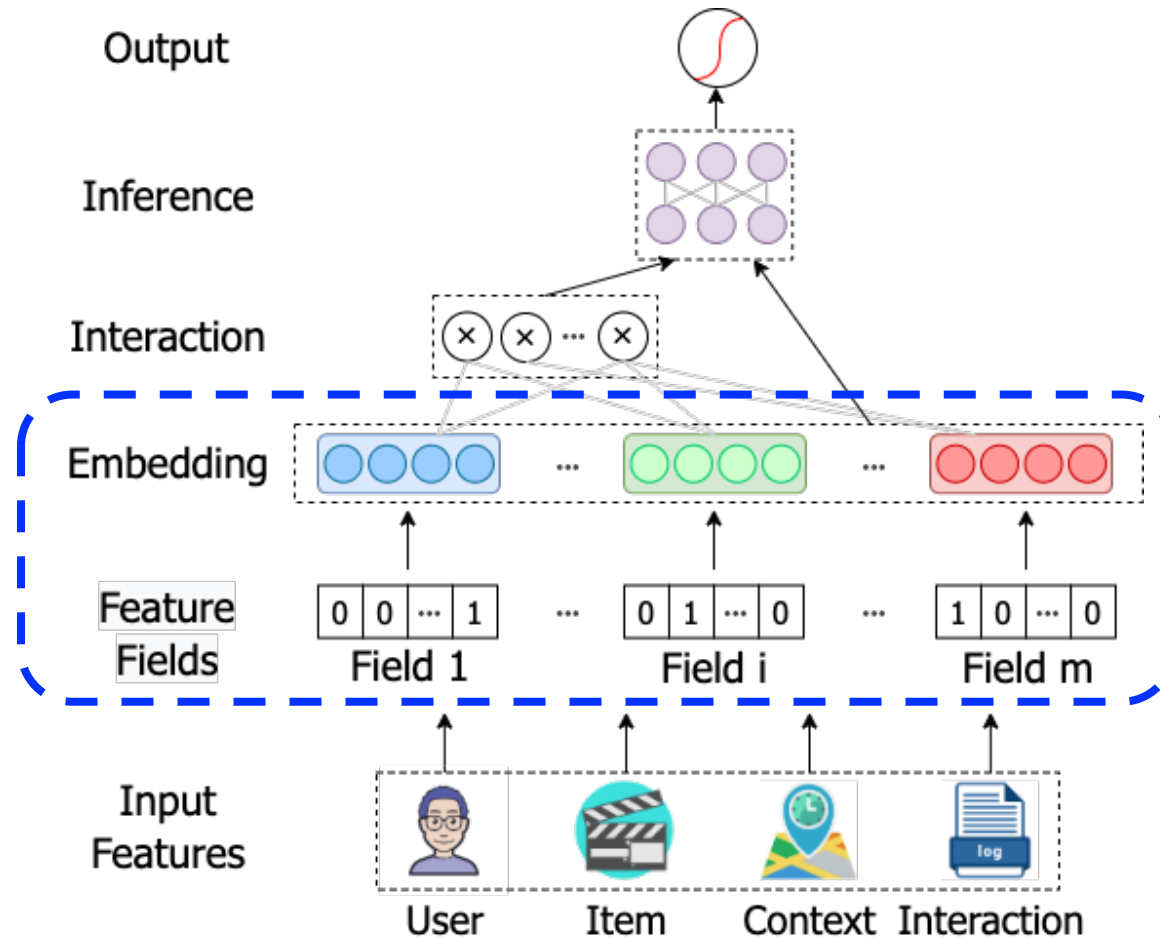
# Neural Architecture Search



- Gradient Descent-based NAS
  - Soft selection on candidate components, weighted sum them
  - Directly update the controller weights by minimizing the loss on validation set



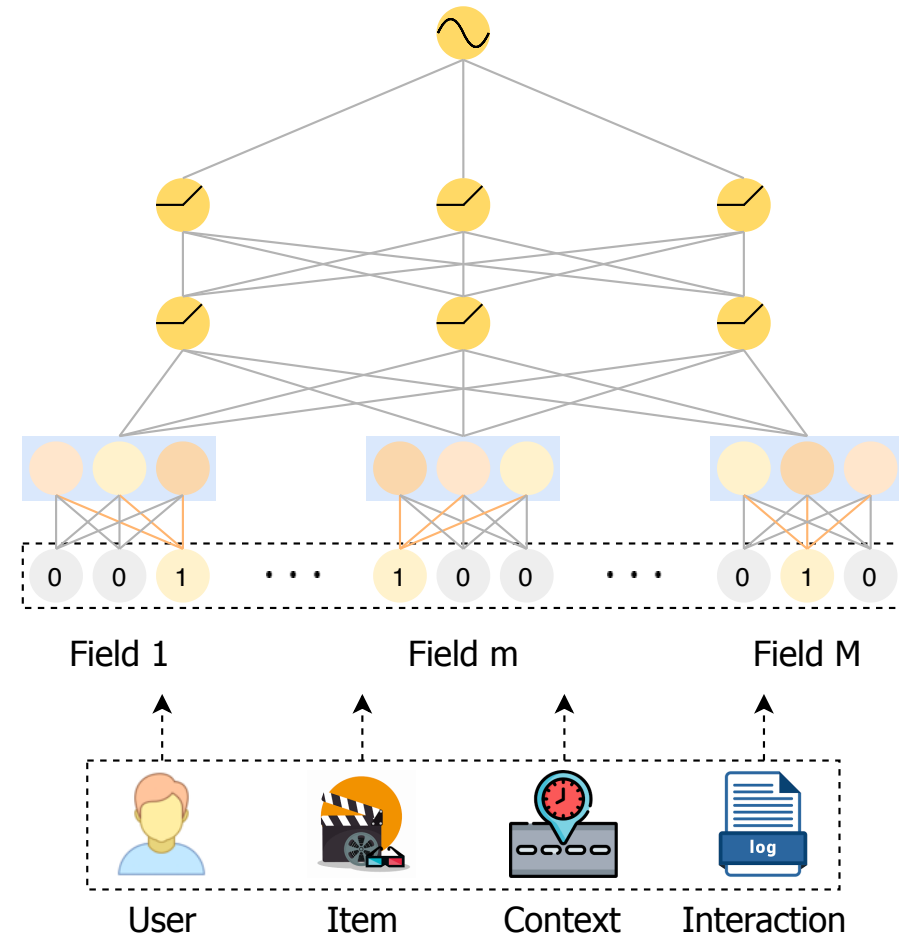
# AutoML in Embedding Layer



# Embedding Components

- Real-world recommender systems involve numerous feature fields

- Users
  - e.g., gender and age
- Items
  - e.g., category and price
- Contextual information
  - e.g., time and location
- Their interactions
  - e.g., *users' purchased items at location A*

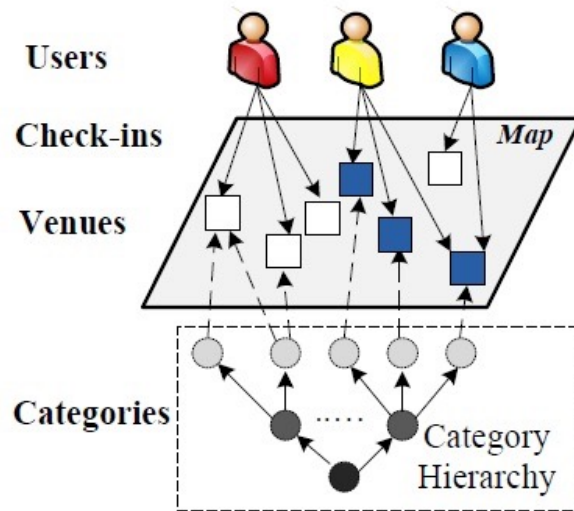
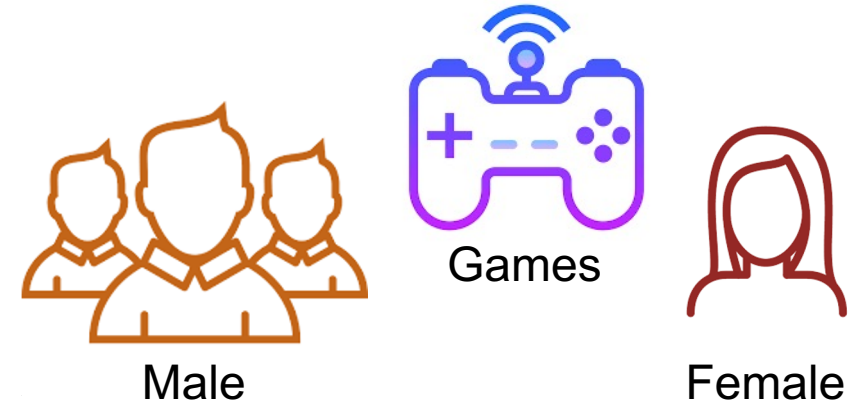


- Features → Embeddings

- Unified** dimension for all features

# Unified Embedding Dimension

- Memory inefficiency problem
  - Embedding dimension  $\rightarrow$  Capacity to encode information
  - Different feature fields have different cardinality
  - Different features have different frequency



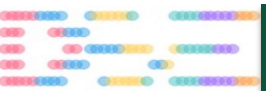
Target	Weekday	Gender	User_ID
1	Tuesday	Male	0000001
0	Monday	Female	3495682
1	Thursday	Female	5676562
0	Friday	Male	9231237

7

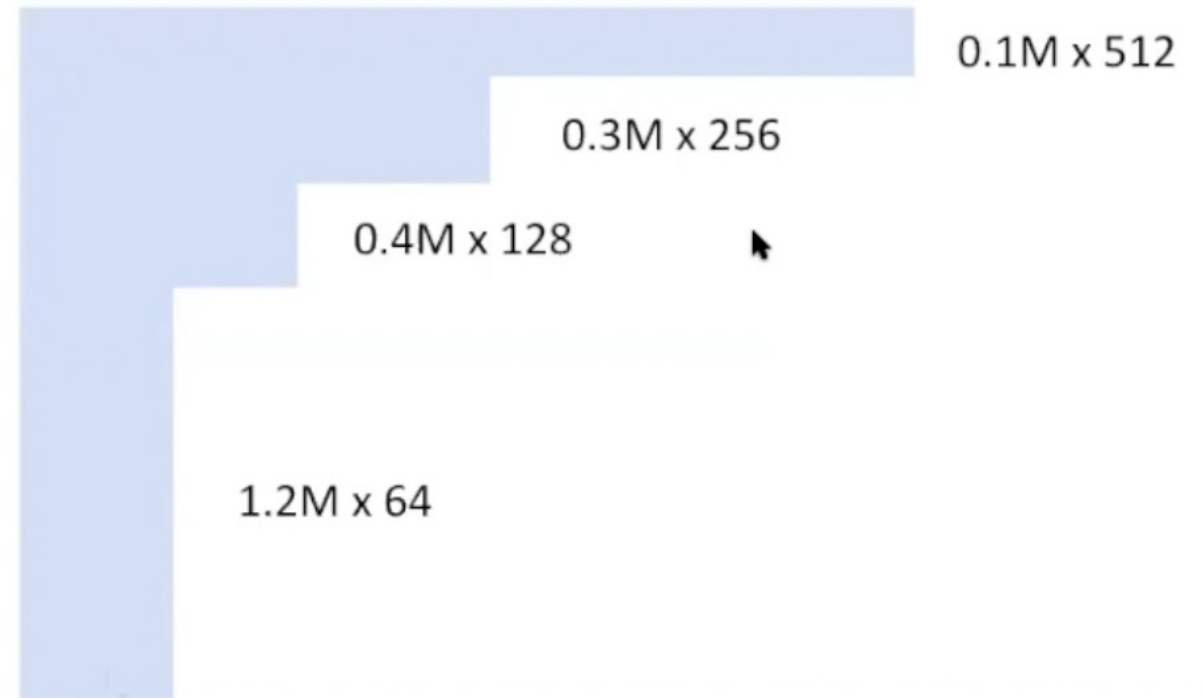
2

million

- **AutoML in Embedding Layer**
  - **NIS** - Neural Input Search for Large Scale Recommendation Models (KDD'2020)
  - **ESAPN** - Automated Embedding Size Search in Deep Recommender Systems (SIGIR'2020)
  - **AutoDim** - Field-aware Embedding Dimension Search in Recommender Systems (WWW'2021)
  - **AutoDis** - Automatic Discretization for Embedding Numerical Features in CTR Prediction (AAAI'2021)
- **AutoML in Interaction Layer**
  - **AutoFIS** - Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction (KDD'2020)
  - **AutoGroup** - Automatic Feature Grouping for Modelling Explicit High-Order Feature Interactions in CTR Prediction (SIGIR'2020)
  - **AutoCTR** - Towards Automated Neural Interaction Discovery for Click-Through Rate Prediction (KDD'2020)



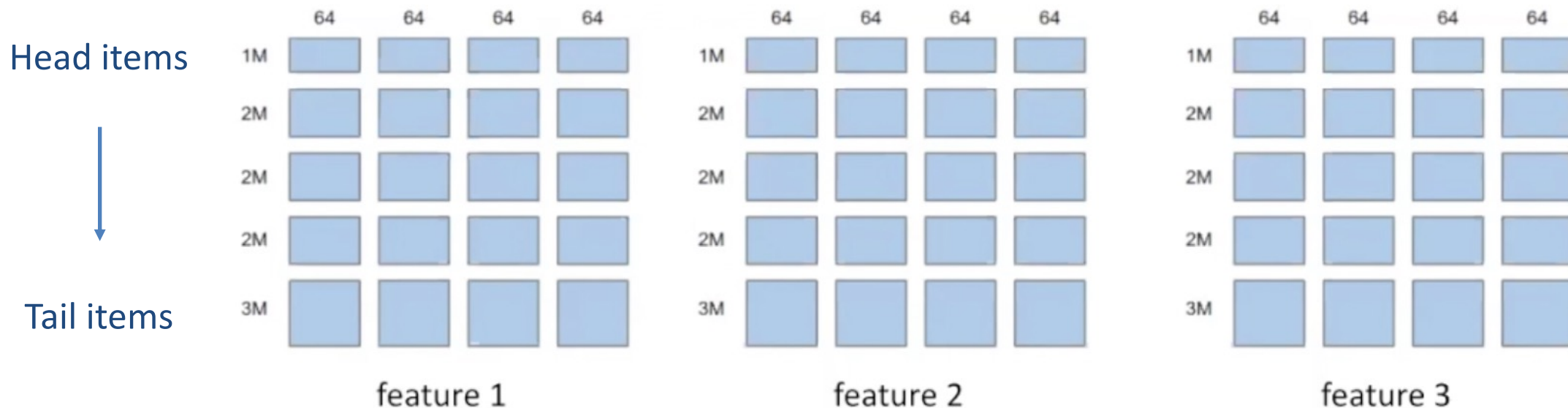
- Head items
  - More data, more information
  - Needing larger embedding size
- Tail items
  - Less data, less information
  - Small embedding size is enough





# NIS - Search Space

- Assume 3 features, each with largest allowed embedding matrix of size 10M x 256
  - Items should be sorted by their frequency
  - Cutting the embedding matrix into smaller pieces
  - The way to cut the embedding matrix is pre-defined

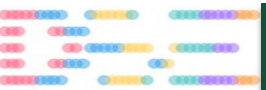
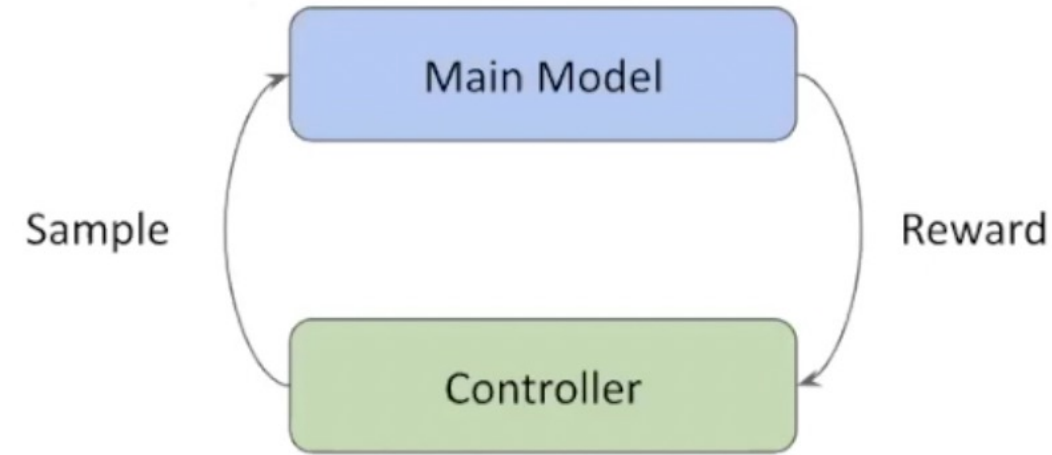




# NIS - Multisize Embedding

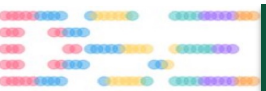
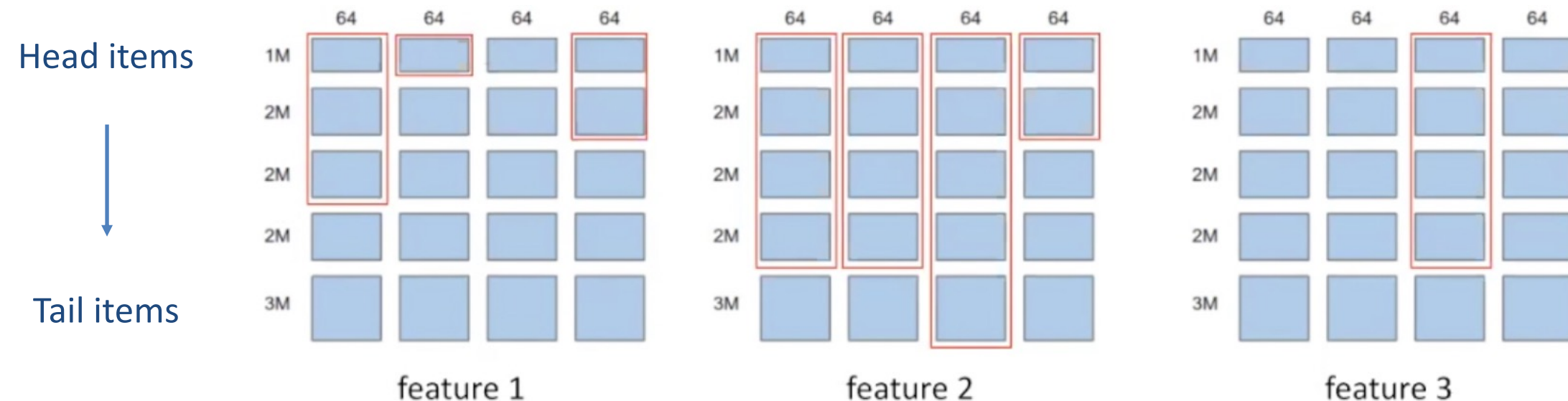
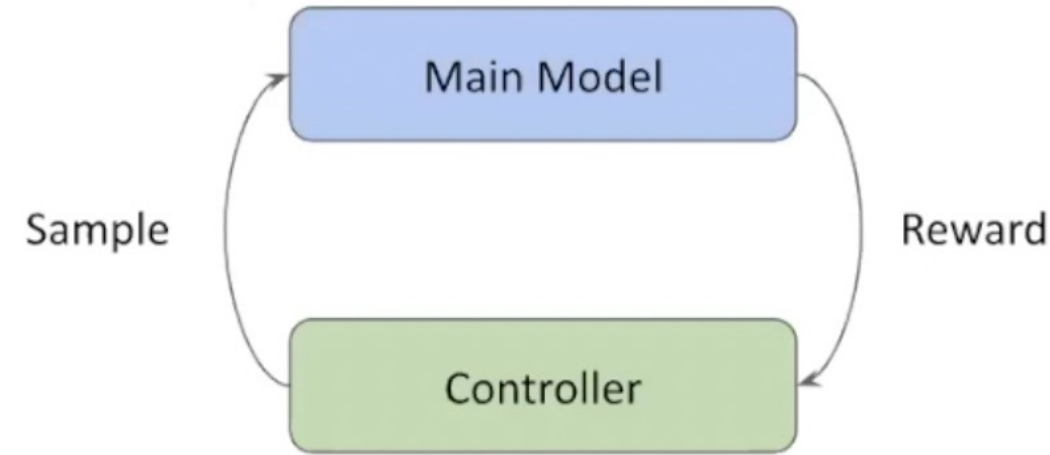


- RL-based AutoML approach
  - Main model is the deep recommendation model
  - Controller learns to sample embedding dimensions that generate higher reward

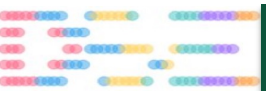


# NIS - Multisize Embedding

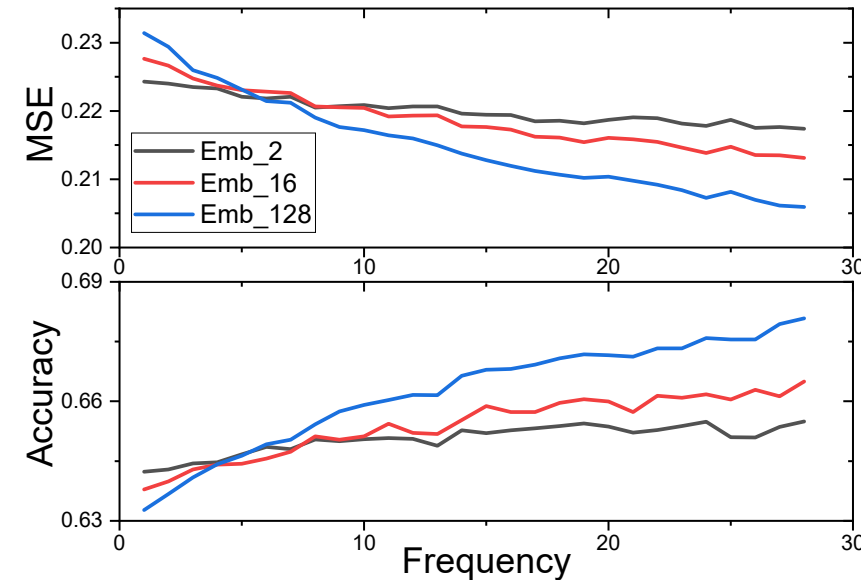
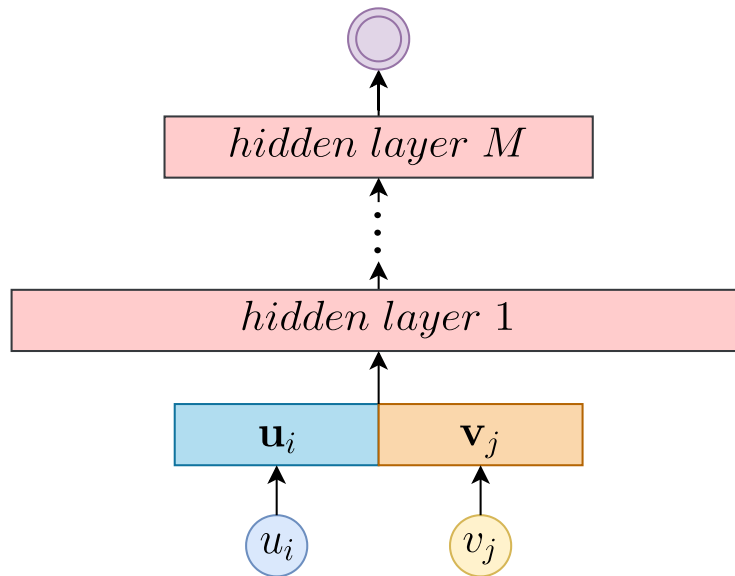
- RL-based AutoML approach
  - Main model is the deep recommendation model
  - Controller learns to sample embedding dimensions that generate higher reward
  - E.g. feature 1: 1M x 192 + 2M x 128 + 2M x 64
  - Reward:**  $R = R_Q - \lambda * C_M$



- **AutoML in Embedding Layer**
  - **NIS** - Neural Input Search for Large Scale Recommendation Models (KDD'2020)
  - **ESAPN** - Automated Embedding Size Search in Deep Recommender Systems (SIGIR'2020)
  - **AutoDim** - Field-aware Embedding Dimension Search in Recommender Systems (WWW'2021)
  - **AutoDis** - Automatic Discretization for Embedding Numerical Features in CTR Prediction (AAAI'2021)
- **AutoML in Interaction Layer**
  - **AutoFIS** - Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction (KDD'2020)
  - **AutoGroup** - Automatic Feature Grouping for Modelling Explicit High-Order Feature Interactions in CTR Prediction (SIGIR'2020)
  - **AutoCTR** - Towards Automated Neural Interaction Discovery for Click-Through Rate Prediction (KDD'2020)



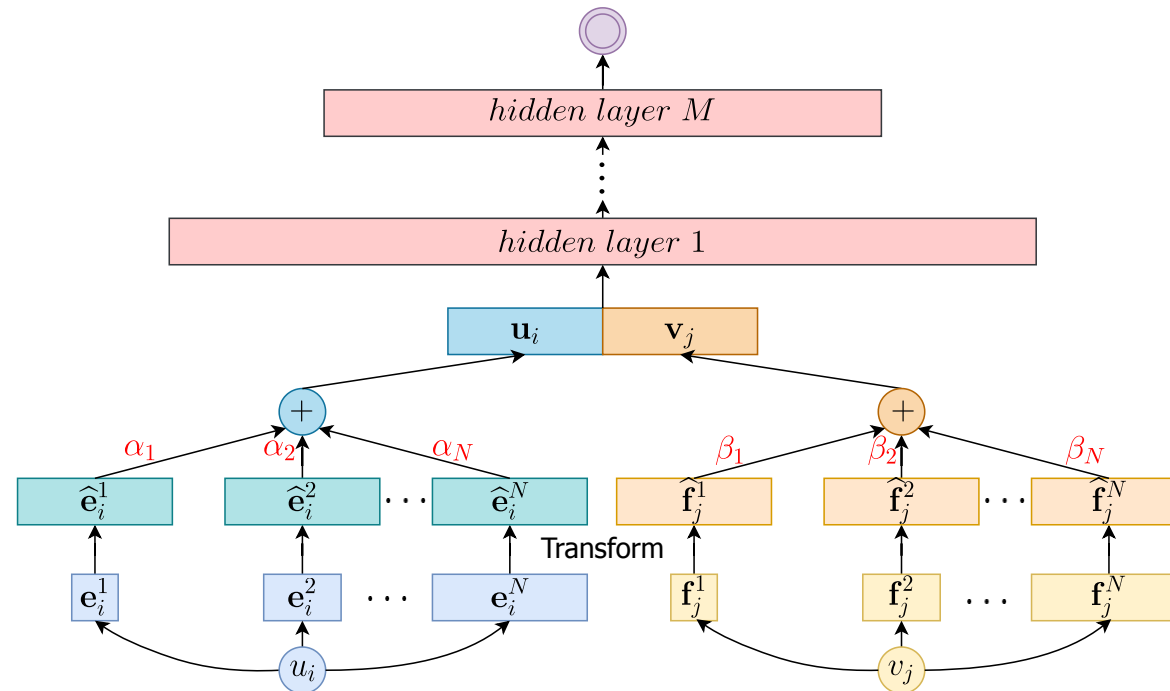
- Preliminary Experiment
  - Frequency: # interactions a user/item



- Embedding dimension often determines the capacity to encode information

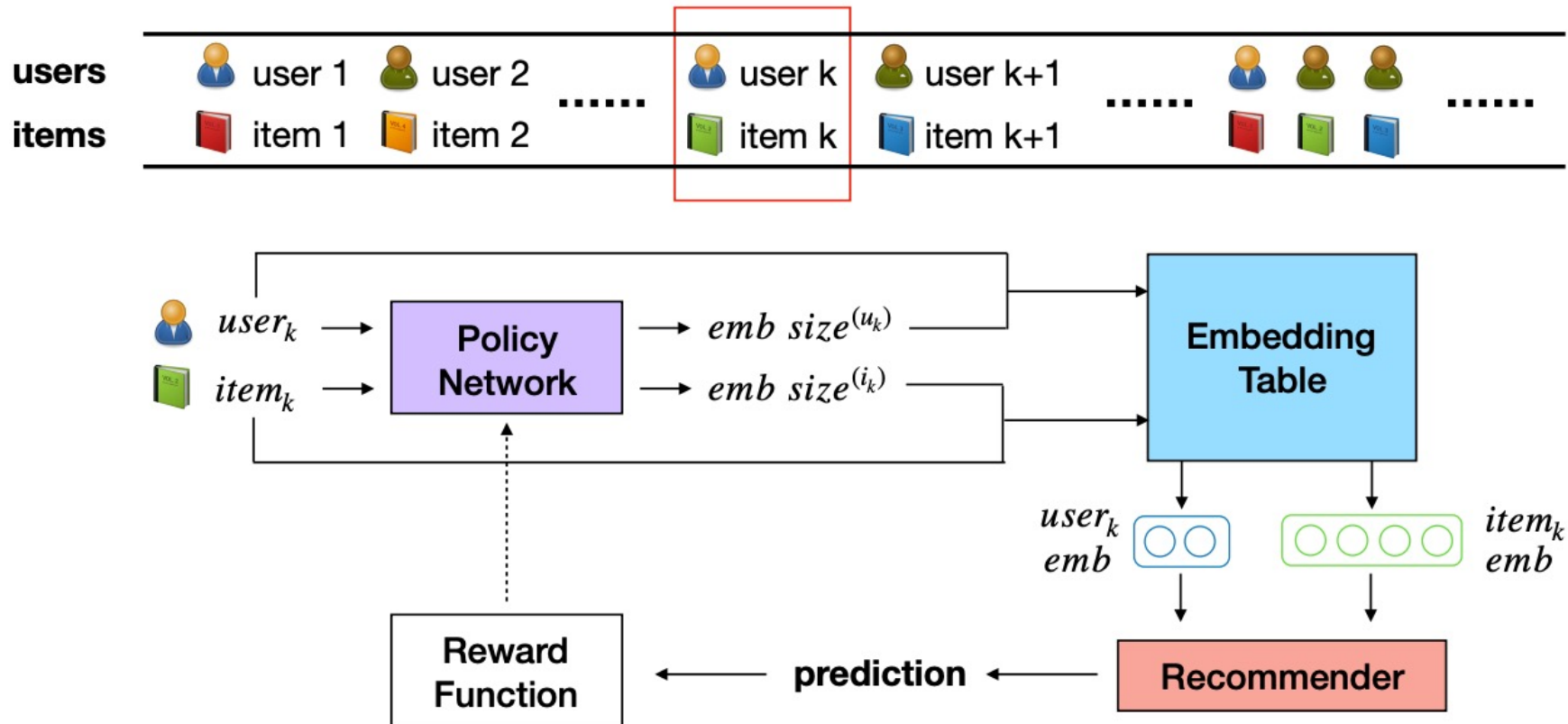
# Motivations

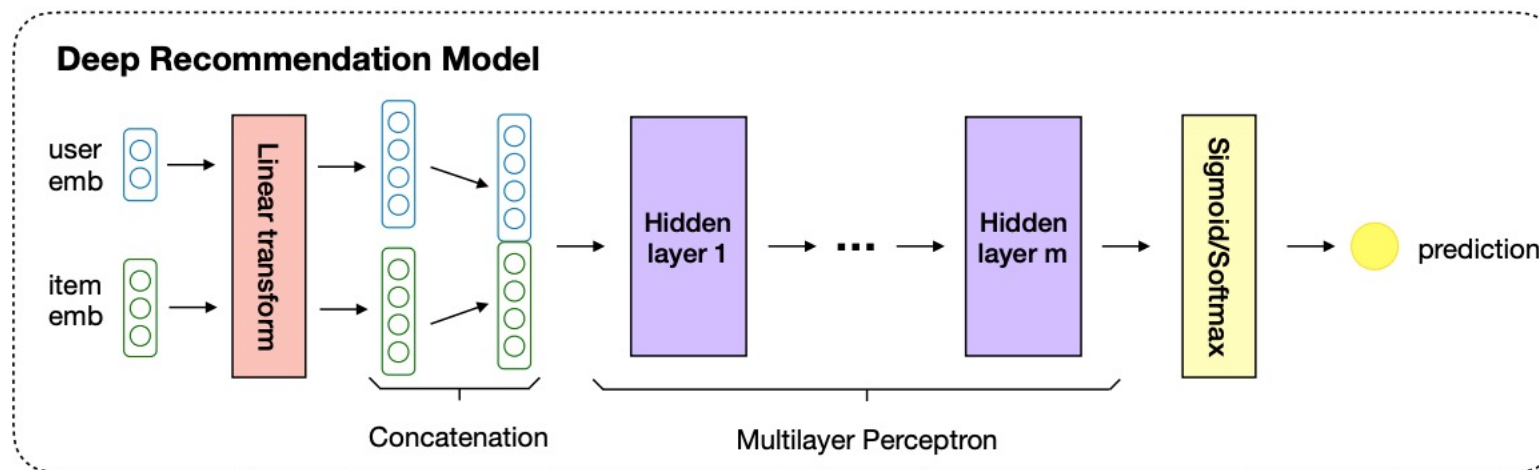
- Dynamically search the embedding sizes for different users and items
  - Optimal recommendation quality all the time
  - More efficient in memory



## Two Components

- Deep recommendation model
- Embedding Size Adjustment Policy Network (ESAPN): hard selection via RL





- Candidate embedding sizes

$$D = \{d_1, d_2, \dots, d_n\} \quad d_1 < d_2 < \dots < d_n$$

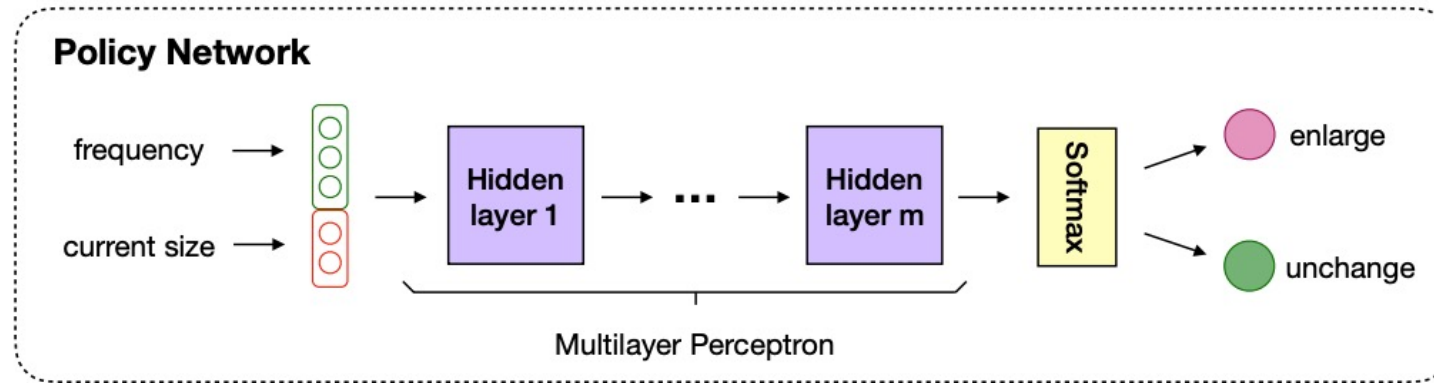
- Linear transformations

$$\mathbf{e}_2 = W_{1 \rightarrow 2} \mathbf{e}_1 + b_{1 \rightarrow 2}$$

$$\mathbf{e}_3 = W_{2 \rightarrow 3} \mathbf{e}_2 + b_{2 \rightarrow 3}$$

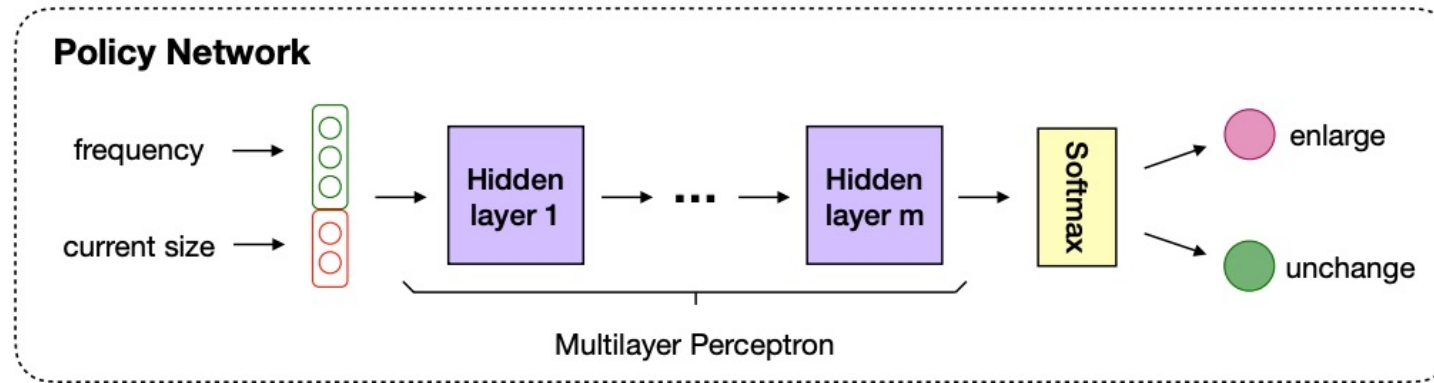
...

$$\mathbf{e}_n = W_{n-1 \rightarrow n} \mathbf{e}_{n-1} + b_{n-1 \rightarrow n}$$



- Environment
  - The deep recommendation model
- State
  - $s = (f, e)$
  - $f$ : frequency       $e$ : current embedding size





- Action
  - Enlarge or Unchange

- Reward

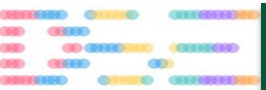
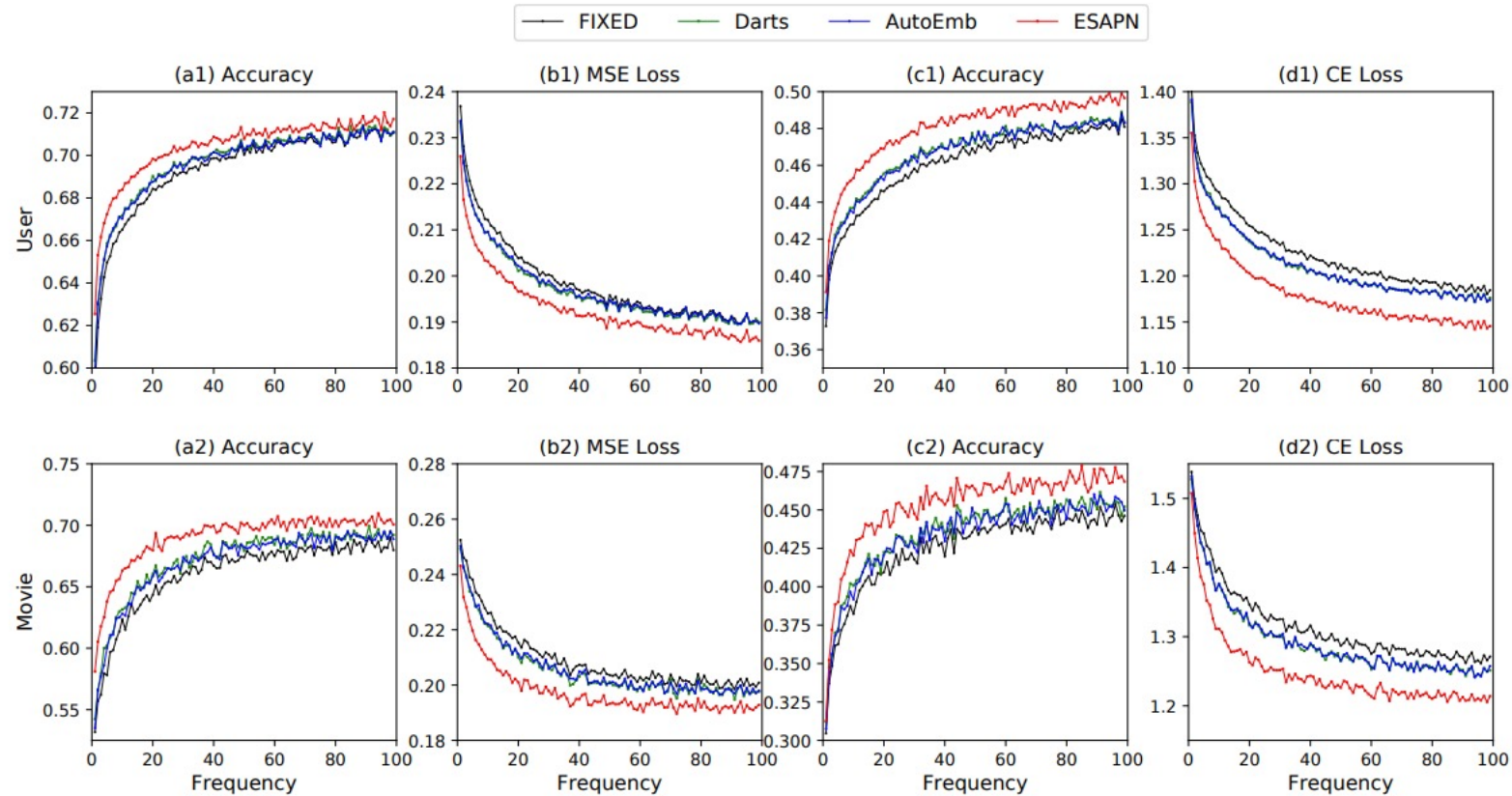
$$L^{(u)} = (L_1^{(u)}, \dots, L_T^{(u)})$$

$$L^{(i)} = (L_1^{(i)}, \dots, L_T^{(i)})$$

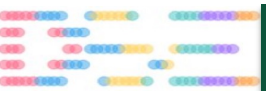
$$R^{(u)} = \frac{1}{T} \sum_{t=1}^T L_t^{(u)} - L$$

$$R^{(i)} = \frac{1}{T} \sum_{t=1}^T L_t^{(i)} - L$$

# Performance with Frequency

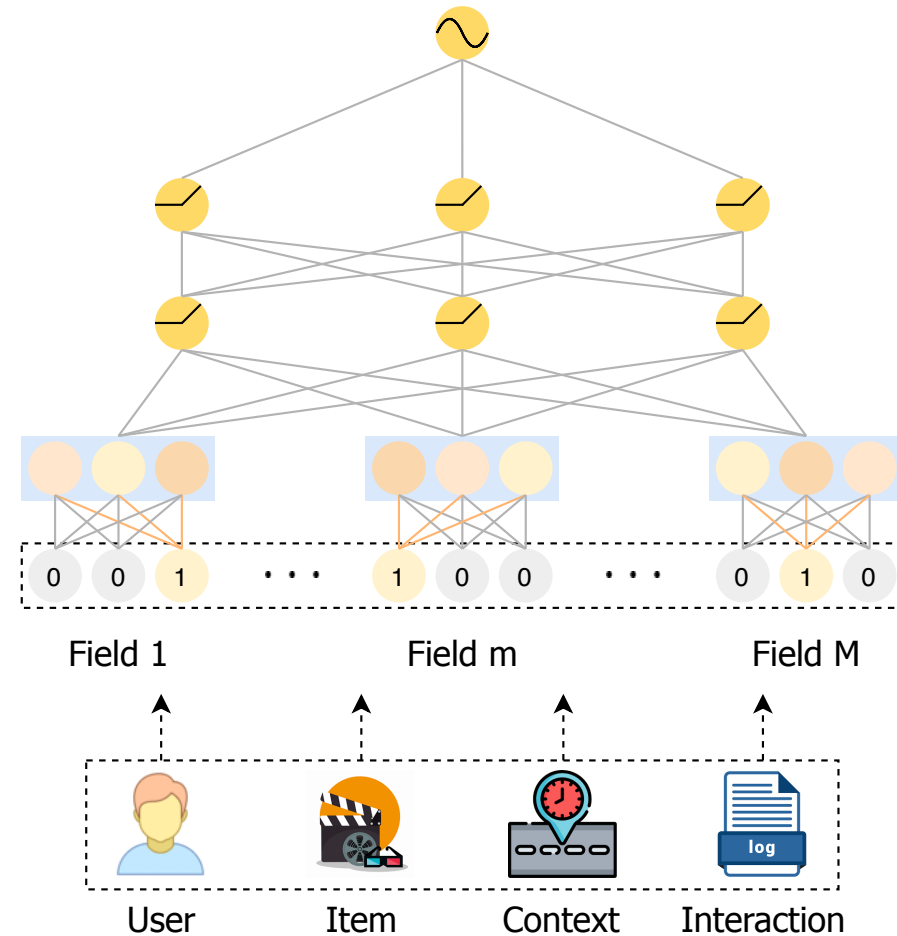


- **AutoML in Embedding Layer**
  - **NIS** - Neural Input Search for Large Scale Recommendation Models (KDD'2020)
  - **ESAPN** - Automated Embedding Size Search in Deep Recommender Systems (SIGIR'2020)
  - **AutoDim** - Field-aware Embedding Dimension Search in Recommender Systems (WWW'2021)
  - **AutoDis** - Automatic Discretization for Embedding Numerical Features in CTR Prediction (AAAI'2021)
- **AutoML in Interaction Layer**
  - **AutoFIS** - Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction (KDD'2020)
  - **AutoGroup** - Automatic Feature Grouping for Modelling Explicit High-Order Feature Interactions in CTR Prediction (SIGIR'2020)
  - **AutoCTR** - Towards Automated Neural Interaction Discovery for Click-Through Rate Prediction (KDD'2020)

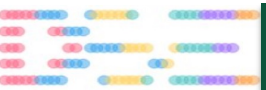
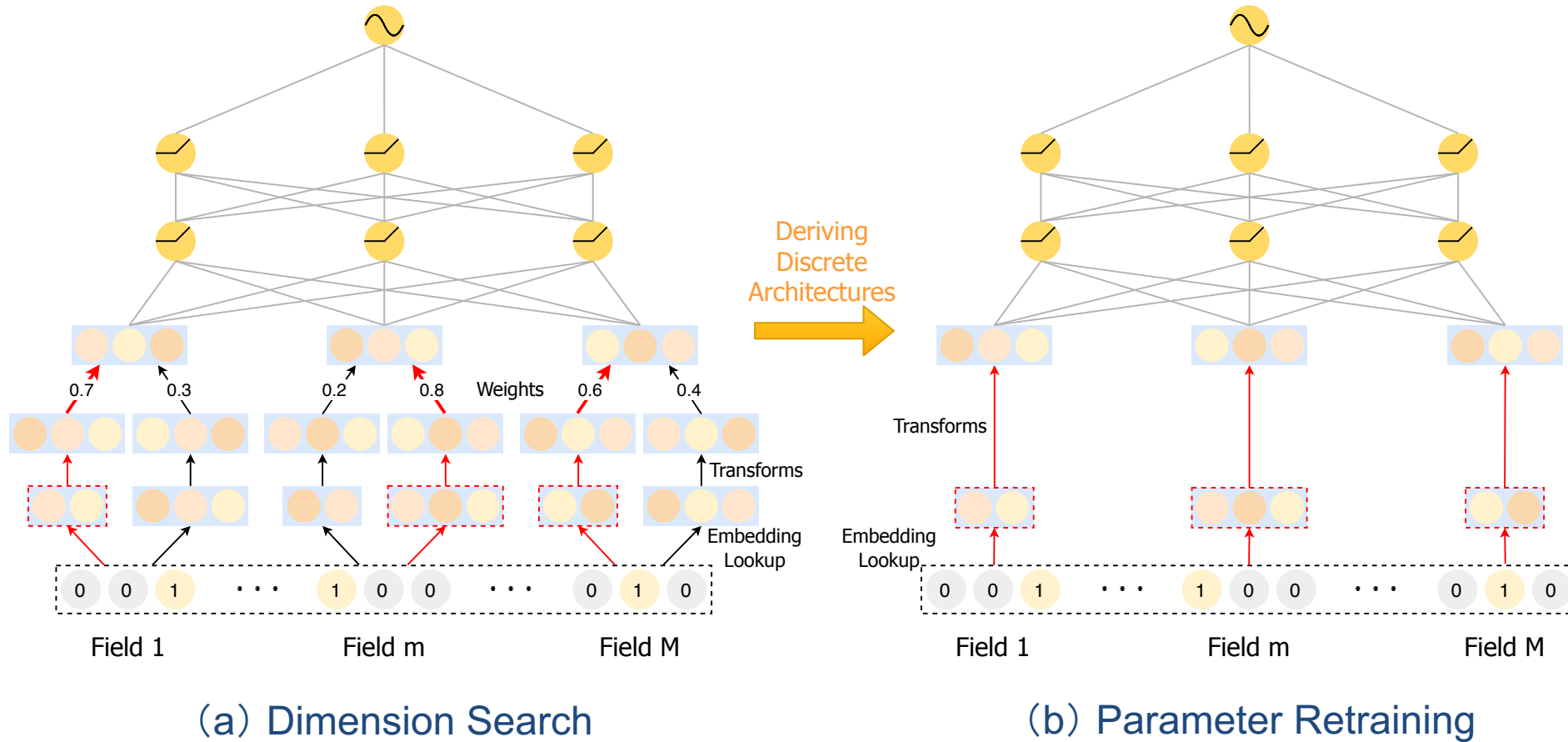


# AutoDim - Motivation

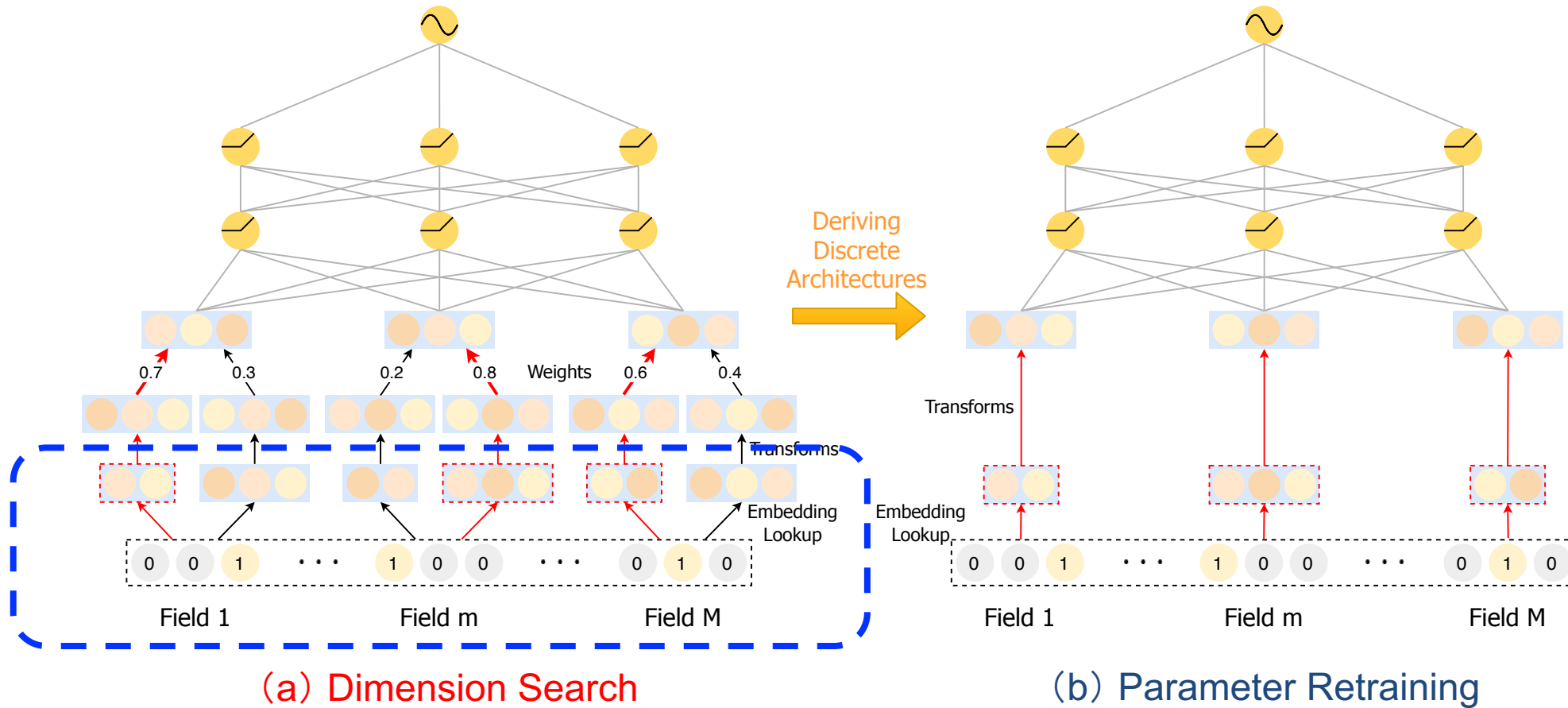
- **Complex** relationship
  - Embedding dimensions
  - Feature distributions
  - Neural network architectures
- **Large** search space
  - M feature field ( $M > 100$ )
  - K candidate dimensions
  - $K^M$  selection space
- **Goal:** Selecting embedding dimensions to different feature fields automatically



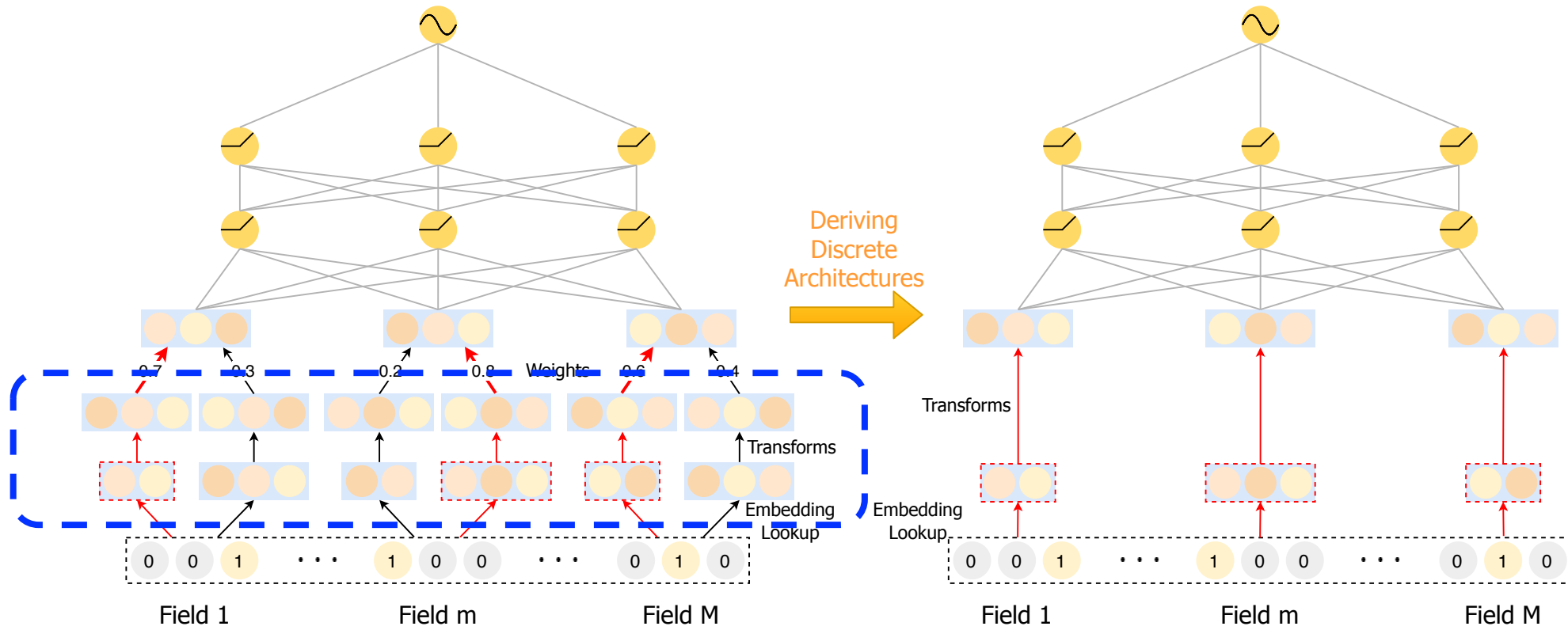
- Two-stage framework



# AutoDim - Dimension Search Stage



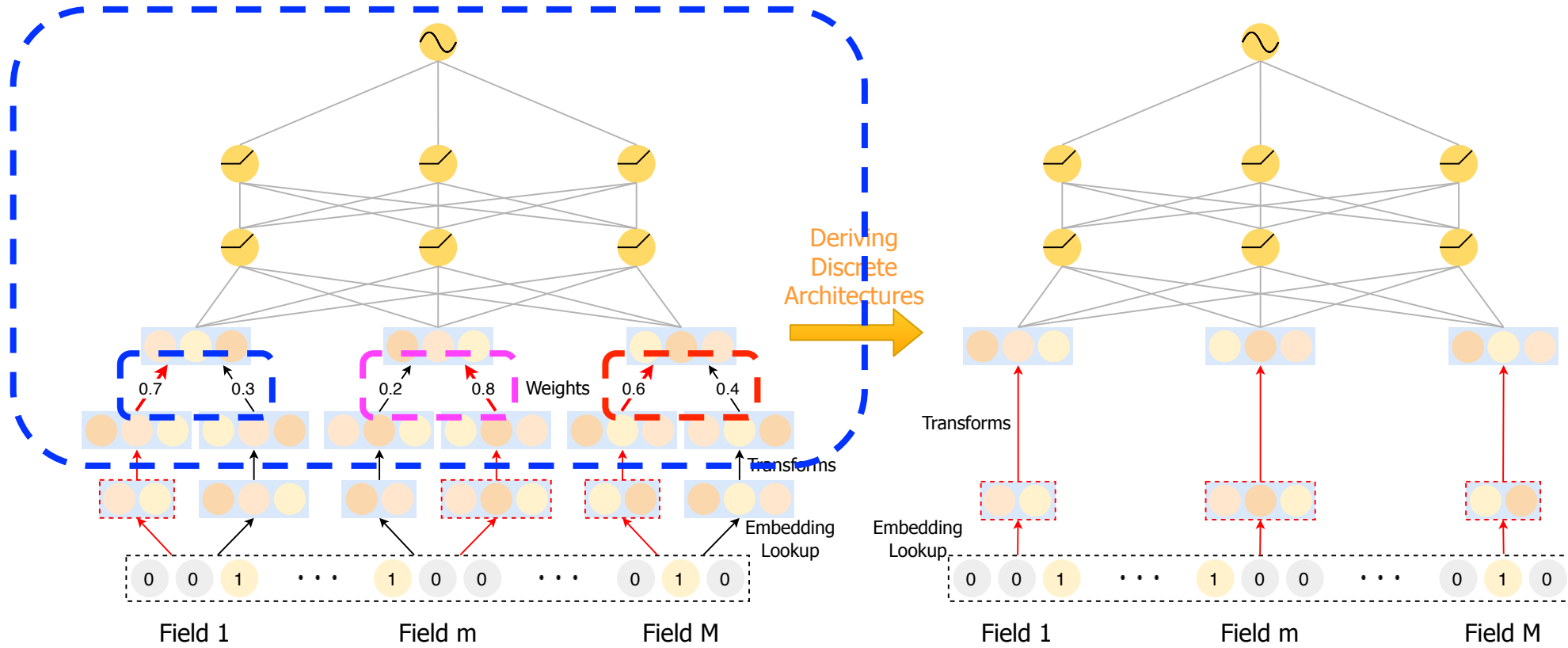
# AutoDim - Dimension Search Stage



(a) Dimension Search

(b) Parameter Retraining

# AutoDim - Dimension Search Stage

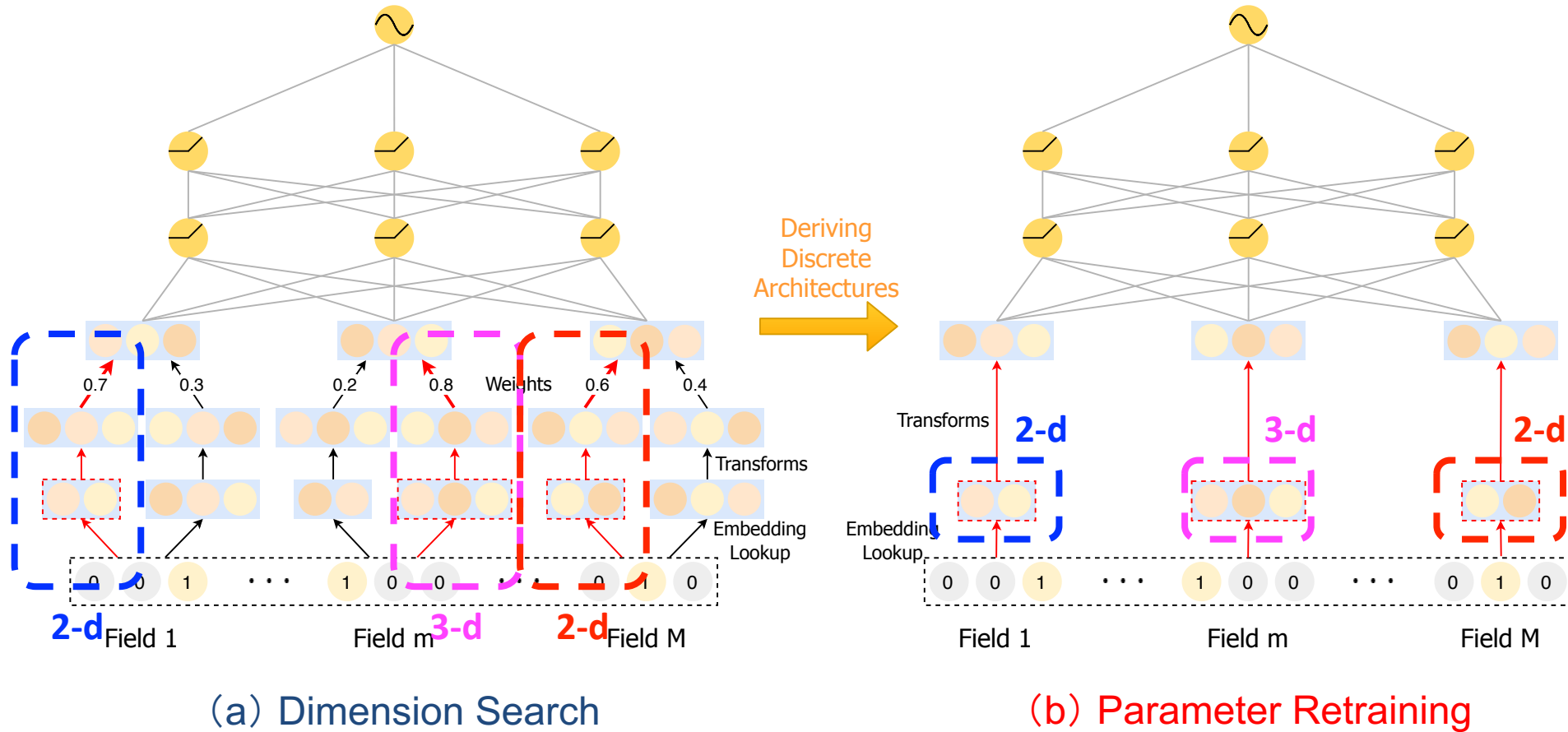


(a) Dimension Search

(b) Parameter Retraining



# AutoDim - Parameter Retraining Stage

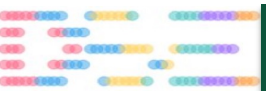


Dataset	Model	Metrics	Search Methods								
			FDE	MDE	DPQ	NIS	MGQE	AEmb	RaS	AD-s	AutoDim
Criteo	FM	AUC	0.8020	0.8027	0.8035	0.8042	0.8046	0.8049	0.8056	0.8063	<b>0.8078*</b>
		Logloss	0.4487	0.4481	0.4472	0.4467	0.4462	0.4460	0.4457	0.4452	<b>0.4438*</b>
		EP (M)	34.778	15.520	20.078	13.636	12.564	13.399	16.236	31.039	<b>11.632*</b>
Criteo	W&D	AUC	0.8045	0.8051	0.8058	0.8067	0.8070	0.8072	0.8076	0.8081	<b>0.8098*</b>
		Logloss	0.4468	0.4464	0.4457	0.4452	0.4446	0.4445	0.4443	0.4439	<b>0.4419*</b>
		EP (M)	34.778	18.562	22.628	14.728	15.741	15.987	18.233	30.330	<b>12.455*</b>
Criteo	DeepFM	AUC	0.8056	0.8060	0.8067	0.8076	0.8080	0.8082	0.8085	0.8089	<b>0.8101*</b>
		Logloss	0.4457	0.4456	0.4449	0.4442	0.4439	0.4438	0.4436	0.4432	<b>0.4416*</b>
		EP (M)	34.778	17.272	25.737	12.955	13.059	13.437	17.816	31.770	<b>11.457*</b>

“\*” indicates the statistically significant improvements (i.e., two-sided t-test with  $p < 0.05$ ) over the best baseline. (M=Million)

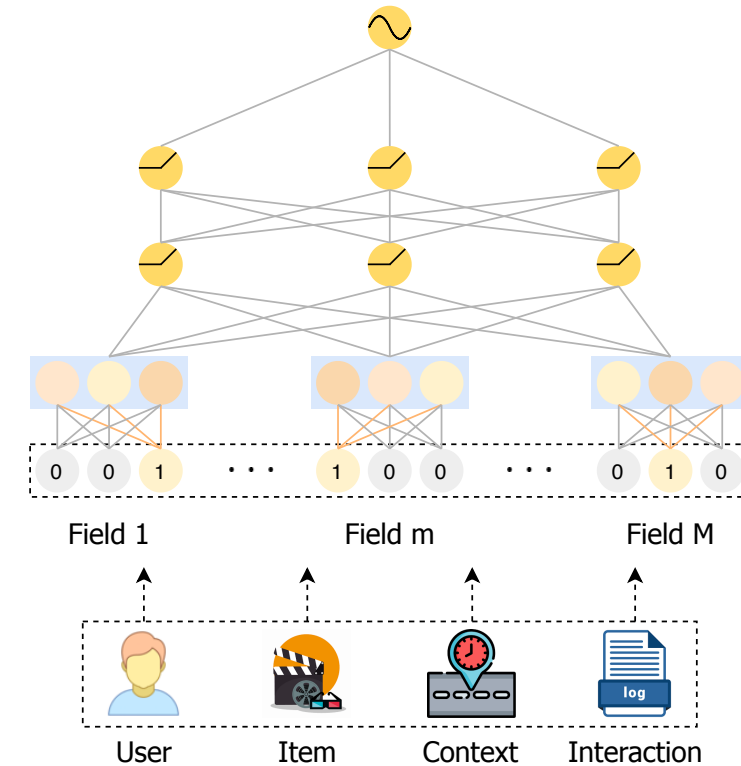
- Metrics: AUC  $\uparrow$ , Logloss  $\downarrow$ , EP  $\downarrow$  (embedding parameters)
- AutoDim is general for **any** deep recommender systems with embedding layer
- **Small** search space: 5 candidate for each feature field
- **AutoDim**  $\rightarrow$  Best AUC and Logloss, and **saving 70~80% embedding parameters**

- **AutoML in Embedding Layer**
  - **NIS** - Neural Input Search for Large Scale Recommendation Models (KDD'2020)
  - **ESAPN** - Automated Embedding Size Search in Deep Recommender Systems (SIGIR'2020)
  - **AutoDim** - Field-aware Embedding Dimension Search in Recommender Systems (WWW'2021)
  - **AutoDis** - Automatic Discretization for Embedding Numerical Features in CTR Prediction (AAAI'2021)
- **AutoML in Interaction Layer**
  - **AutoFIS** - Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction (KDD'2020)
  - **AutoGroup** - Automatic Feature Grouping for Modelling Explicit High-Order Feature Interactions in CTR Prediction (SIGIR'2020)
  - **AutoCTR** - Towards Automated Neural Interaction Discovery for Click-Through Rate Prediction (KDD'2020)



- Real-world recommender systems involve numerous feature fields

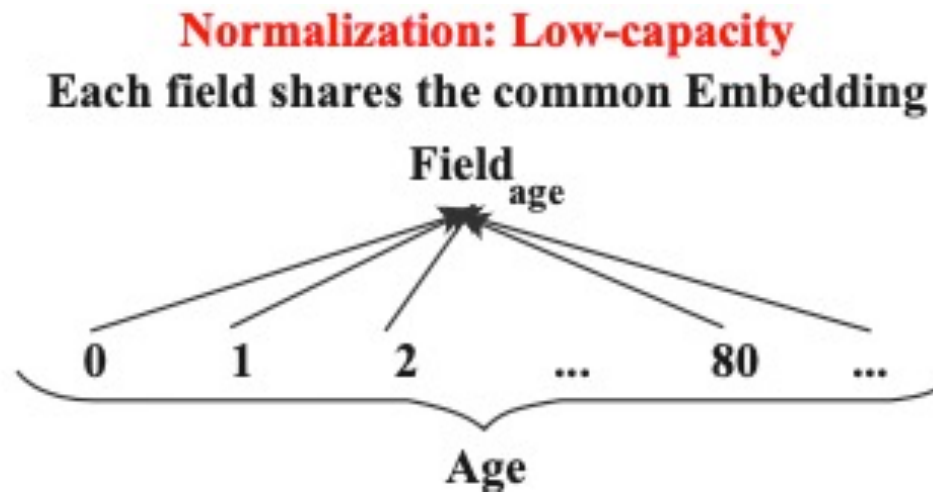
- Users
  - e.g., gender and age
- Items
  - e.g., category and price
- Contextual information
  - e.g., time and location
- Their interactions
  - e.g., *users'* purchased *items* at *location A*



- E.g., Gender=Male, Day=Tuesday, Height=175.6, Age=18

- Categorical field Gender v.s. Numerical field Height

- Normalization
  - All the numerical features in the same field share a single embedding and scalar multiply with their values



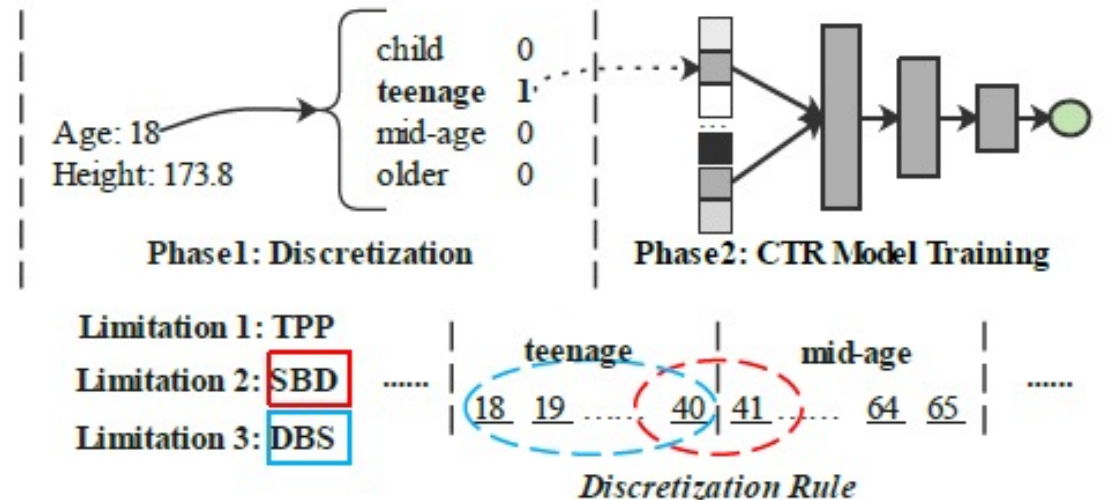
- Disadvantage
  - Assuming embeddings of different features in the same field are linearly related to each other

# Existing Methods for Numerical features



## Discretization

- E.g., partitioning the range of the feature values into  $k$  buckets

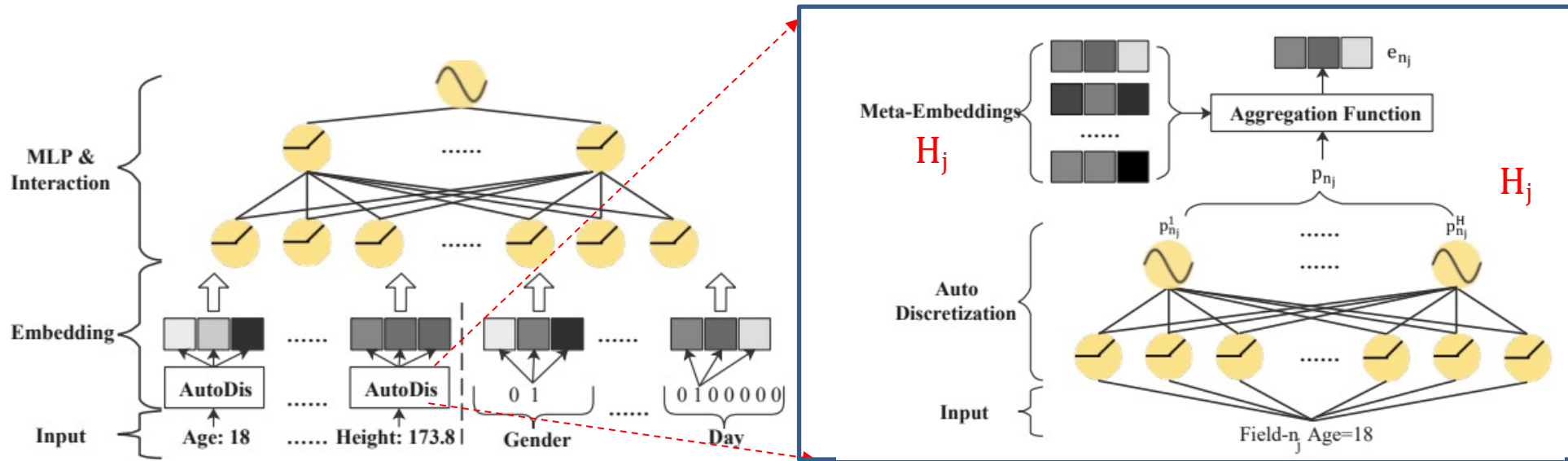


## Disadvantages

- TPP: cannot be optimized together with main model
- SBD: different embeddings for similar numerical value 40 and 41
- DBS: same embeddings for very different numerical value 18 and 40



# Aggregation Function



## Automatic Discretization

Continuous value mapping :

$$\hat{x}_{n_j}^h = \mathbf{W}_{n_j}^h \cdot x_{n_j}$$

Softmax :

$$p_{n_j}^h = \frac{e^{\frac{1}{\tau} \hat{x}_{n_j}^h}}{\sum_{l=1}^{H_j} e^{\frac{1}{\tau} \hat{x}_{n_j}^l}}$$

Soft discretization output:

$$g(x_{n_j}) = [p_{n_j}^1, \dots, p_{n_j}^h, \dots, p_{n_j}^{H_j}]$$

## Aggregation Function

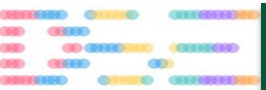
Max-Pooling:  $e_{n_j} = \mathbf{ME}_{n_j}^k$ , where  $k = \arg \max_{h \in \{1, 2, \dots, H_j\}} p_{n_j}^h$ ,

Top-K-Sum:  $e_{n_j} = \sum_{l=1}^K \mathbf{ME}_{n_j}^{k_l}$ , where  $k_l = \arg \text{top}_l_{h \in \{1, 2, \dots, H_j\}} p_{n_j}^h$ ,

Weighted-Average:  $e_{n_j} = \sum_{l=1}^{H_j} p_{n_j}^l \cdot \mathbf{ME}_{n_j}^l$ .

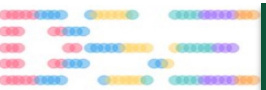
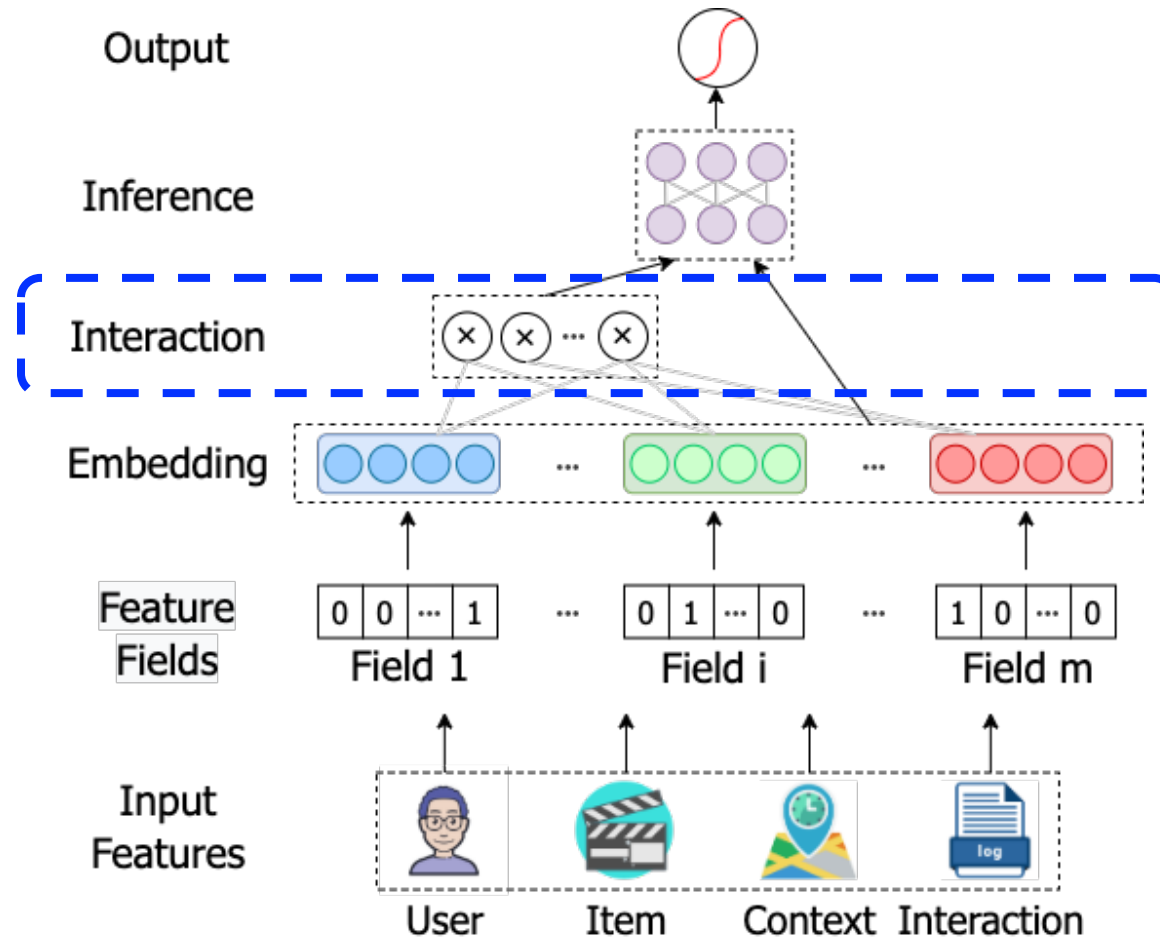


- End-To-End:
  - The discretization of numerical features can be optimized jointly with the main model
- Continuous-But-Different
  - Different feature values are assigned with different embeddings
  - Closer the feature values have more similar the embeddings





# AutoML in Interaction Layer



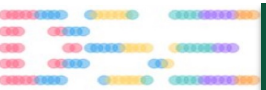
# Background



- Multi-field data

Target	Weekday	Gender	City	Product Category
1	Tuesday	Male	London	Sports
0	Monday	Female	New York	Cosmetics
1	Thursday	Female	Beijing	Clothing
0	Friday	Male	Tokyo	Food

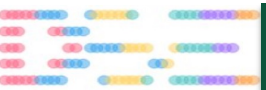
- High dimensional and sparse



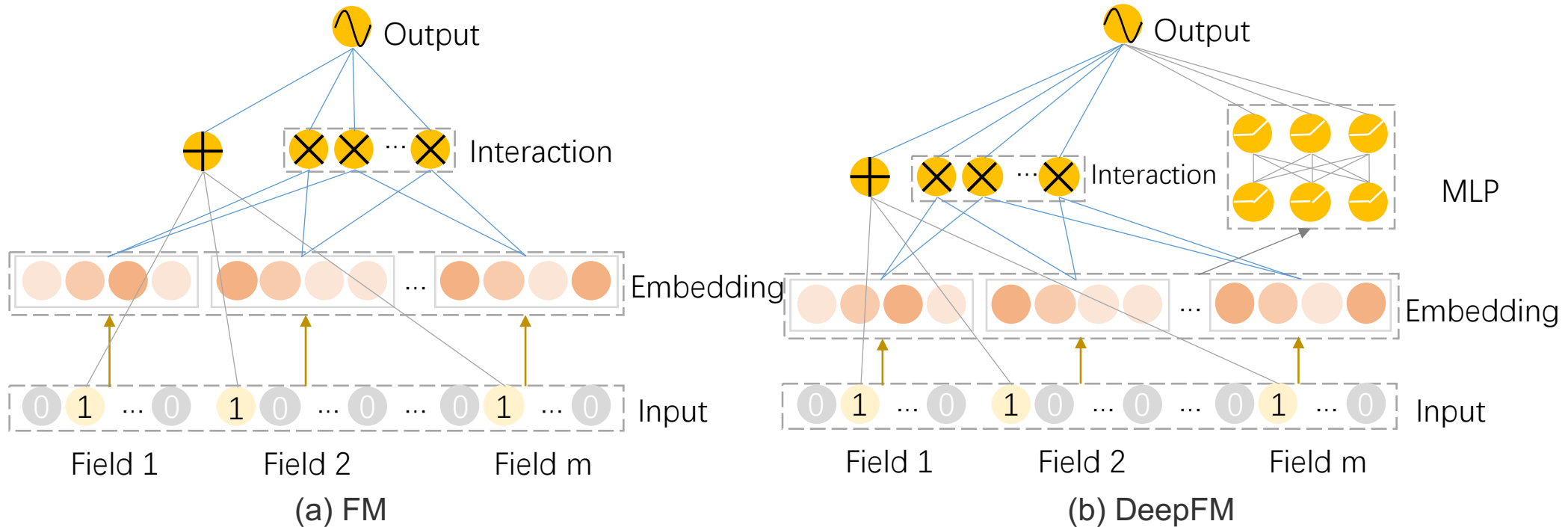
# Effectively Modelling Feature Interactions Is Important



- User behavior is complicated to model
- Both low-order and high-order feature interactions play important roles to model user behavior.
  - People like to download popular apps → id of an app may be a signal
  - People often download apps for food delivery at meal time → interaction between app category and time-stamp may be a signal
  - Male teenagers like shooting game or RPG → interaction of app category, user gender and age may be a signal
- Most feature interactions are hidden in data and difficult to identify (e.g., “diaper and beer” rule)



- Factorization models are the models where the interaction of several embeddings from different features is modeled into a real number by some operation such as inner product or neural network

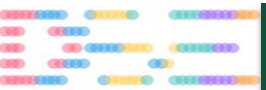


- Enumerate all feature interactions
  - Large memory and computation cost and difficult to be extended into high-order interactions
  - Useless interaction
- Require human efforts to identify important feature interactions
  - high labor cost
  - risks missing some counterintuitive (but important) interactions

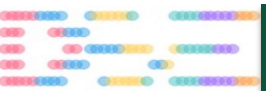
$$y_{\text{FM}}(x) = \text{sigmoid} \left( \sum_{i=1}^N \omega_i x_i + \sum_{i=1}^N \sum_{j=i+1}^N \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j \right)$$

$$y_{\text{FFM}}(x) = \text{sigmoid} \left( \sum_{i=1}^N \omega_i x_i + \sum_{i=1}^N \sum_{j=i+1}^N \langle \mathbf{v}_{i,f(j)}, \mathbf{v}_{j,f(i)} \rangle x_i x_j \right)$$

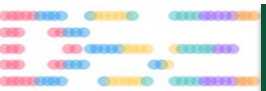
$$y_{\text{FwFM}}(x) = \text{sigmoid} \left( \sum_{i=1}^N \omega_i x_i + \sum_{i=1}^N \sum_{j=i+1}^N A_{f(i),f(j)} \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j \right)$$



- **AutoML in Embedding Layer**
  - **NIS** - Neural Input Search for Large Scale Recommendation Models (KDD'2020)
  - **ESAPN** - Automated Embedding Size Search in Deep Recommender Systems (SIGIR'2020)
  - **AutoDim** - Field-aware Embedding Dimension Search in Recommender Systems (WWW'2021)
  - **AutoDis** - Automatic Discretization for Embedding Numerical Features in CTR Prediction (AAAI'2021)
- **AutoML in Interaction Layer**
  - **AutoFIS** - Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction (KDD'2020)
  - **AutoGroup** - Automatic Feature Grouping for Modelling Explicit High-Order Feature Interactions in CTR Prediction (SIGIR'2020)
  - **AutoCTR** - Towards Automated Neural Interaction Discovery for Click-Through Rate Prediction (KDD'2020)

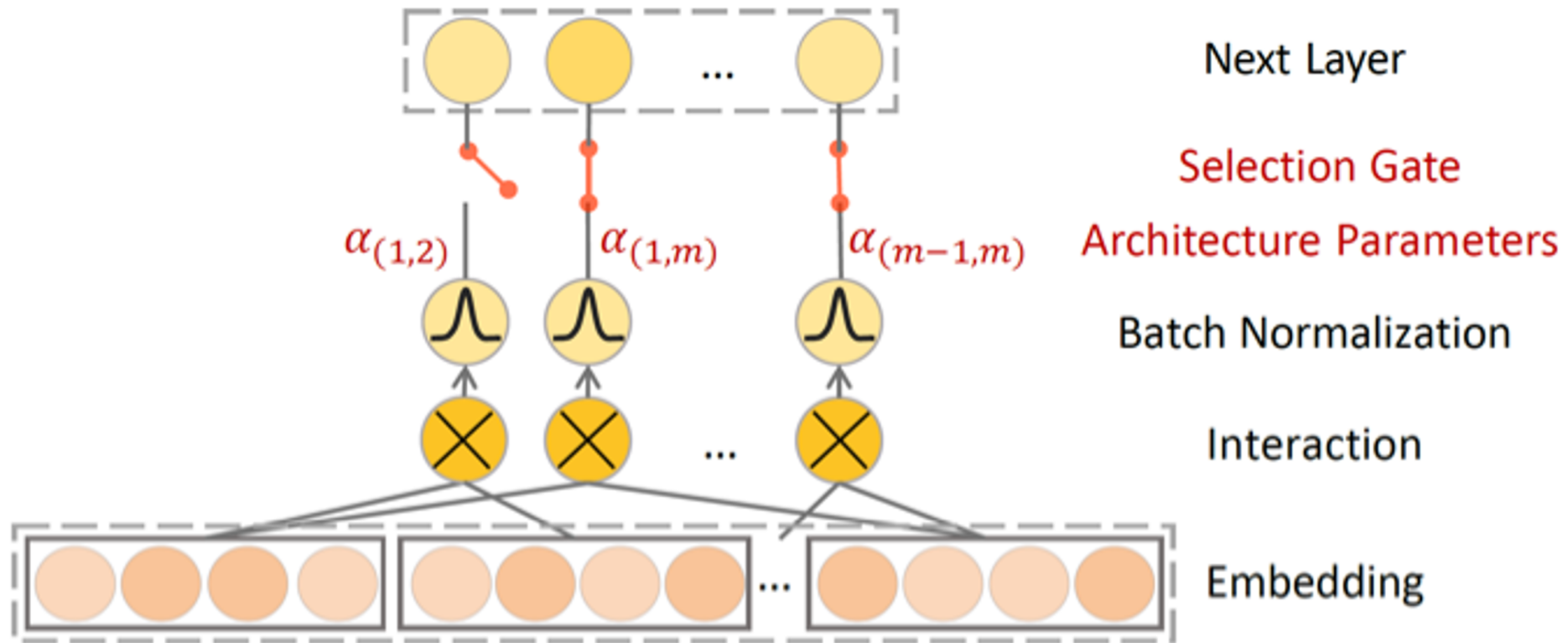


- Search Stage
  - Detect useful feature interactions
- Retrain Stage
  - Retrain model with selected feature interactions



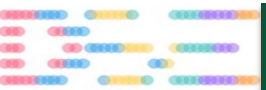
$$l_{\text{AutoFIS}} = \langle w, x \rangle + \sum_{i=1}^m \sum_{j>i}^m \alpha_{(i,j)} \langle e_i, e_j \rangle$$

Indicator  $\alpha = 0$  or  $1$

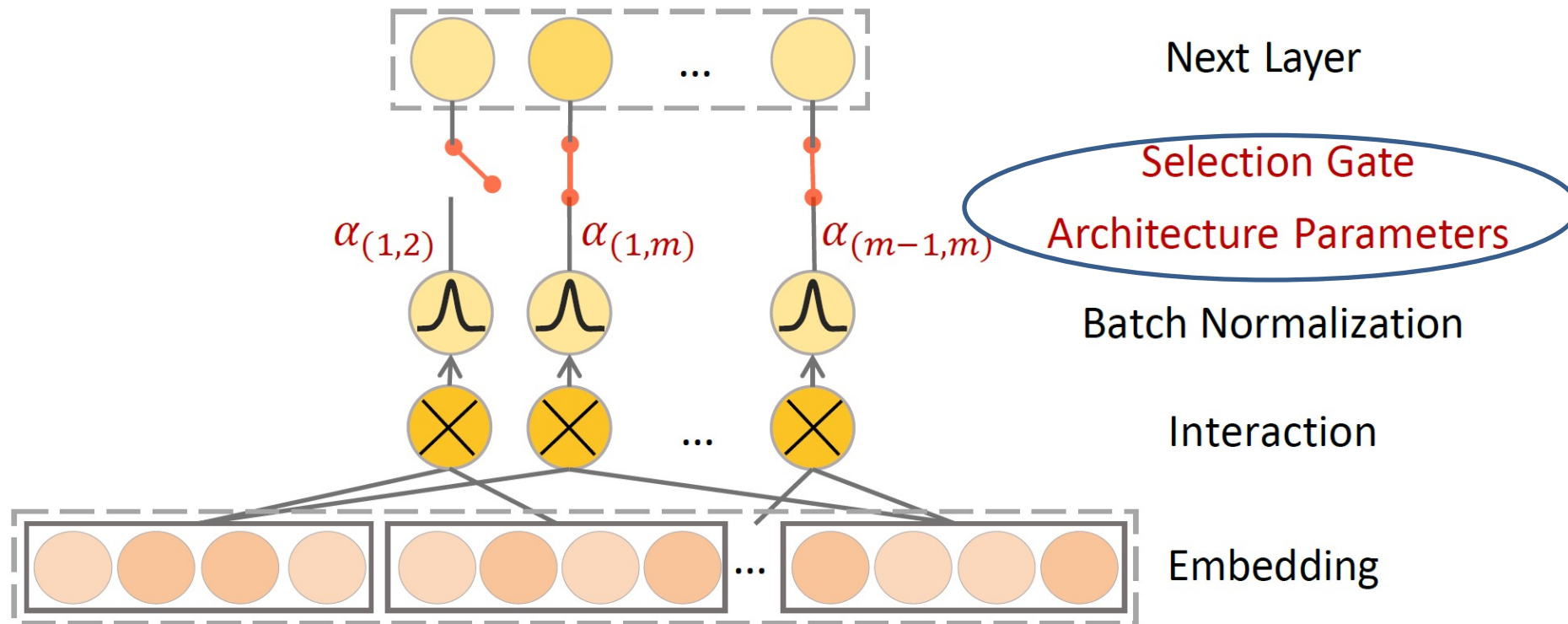




- Gate for each feature interaction
  - Huge search space  $2^{C_m^2}$
  
- Discrete search space -> Continuous search space
  - Architecture parameters  $\alpha$



$$l_{\text{AutoFIS}} = \langle w, x \rangle + \sum_{i=1}^m \sum_{j>i}^m \alpha_{(i,j)} \langle e_i, e_j \rangle$$

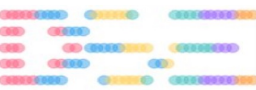


# Experiment Results in Huawei Dataset

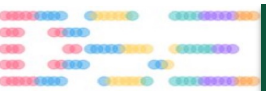


Model	AUC	log loss	top	Rel. Impr
FM	0.8880	0.08881	100%	0
FwFM	0.8897	0.08826	100%	0.19%
AFM	0.8915	0.08772	100%	0.39%
FFM	0.8921	0.08816	100%	0.46%
DeepFM	0.8948	0.08735	100%	0.77%
AutoFM(2nd)	0.8944*	0.08665*	37%	0.72%
AutoDeepFM(2nd)	<b>0.8979*</b>	<b>0.08560*</b>	<b>15%</b>	1.11%

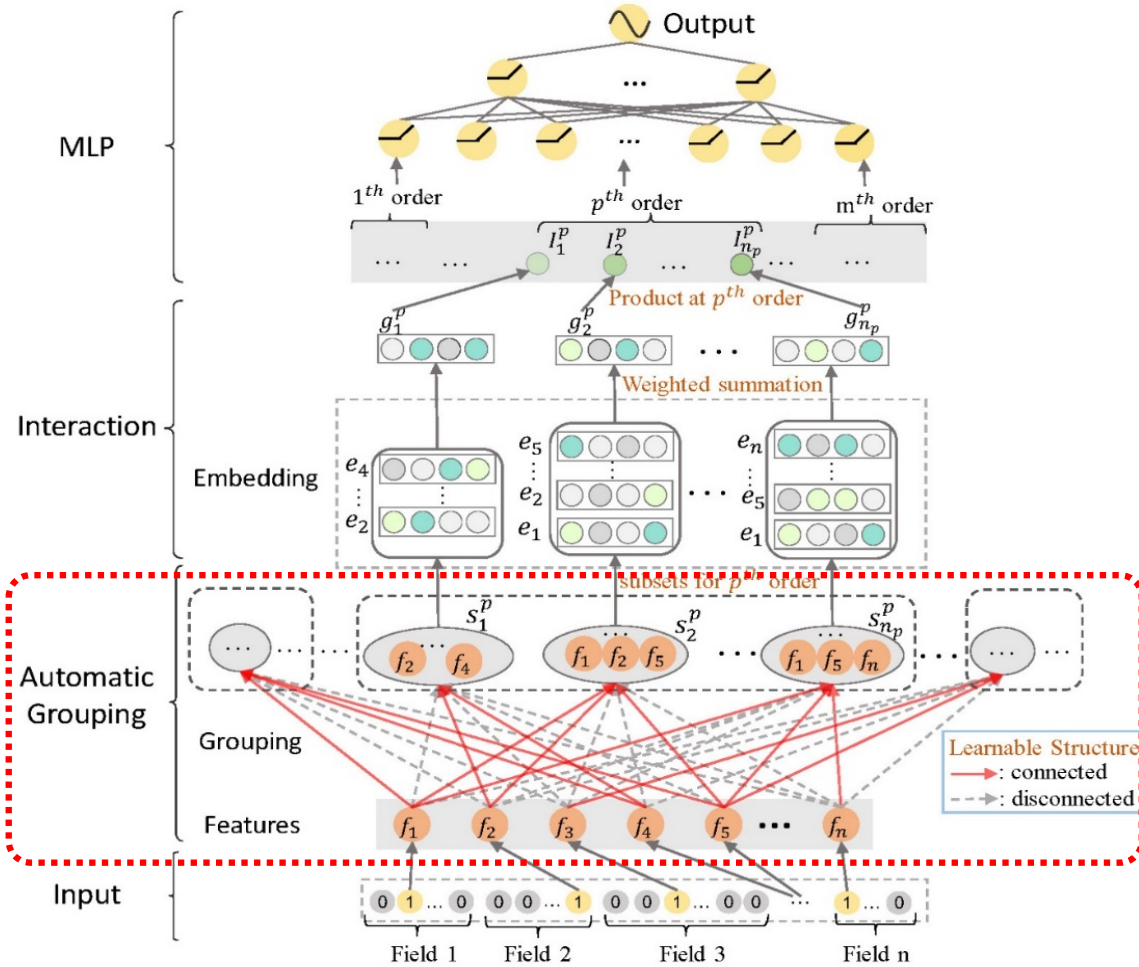
\* denotes statistically significant improvement (measured by t-test with p-value < 0.005). AutoFM compares with FM and AutoDeepFM compares with all baselines.



- **AutoML in Embedding Layer**
  - **NIS** - Neural Input Search for Large Scale Recommendation Models (KDD'2020)
  - **ESAPN** - Automated Embedding Size Search in Deep Recommender Systems (SIGIR'2020)
  - **AutoDim** - Field-aware Embedding Dimension Search in Recommender Systems (WWW'2021)
  - **AutoDis** - Automatic Discretization for Embedding Numerical Features in CTR Prediction (AAAI'2021)
- **AutoML in Interaction Layer**
  - **AutoFIS** - Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction (KDD'2020)
  - **AutoGroup** - Automatic Feature Grouping for Modelling Explicit High-Order Feature Interactions in CTR Prediction (SIGIR'2020)
  - **AutoCTR** - Towards Automated Neural Interaction Discovery for Click-Through Rate Prediction (KDD'2020)



# Automatic Feature Grouping Stage



Each feature is possible to be selected into the feature sets of each order.

- $\Pi_{i,j}^p \in \{0,1\}$ : whether select feature  $f_i$  into the  $j^{th}$  set of order- $p$ .

To make the selection differentiable, we relax the binary discrete value to a softmax over the two possibilities:

$$\bar{\Pi}_{i,j}^p = \frac{1}{1 + \exp(-\alpha_{i,j}^p)} \Pi_{i,j}^p + \frac{\exp(-\alpha_{i,j}^p)}{1 + \exp(-\alpha_{i,j}^p)} (1 - \Pi_{i,j}^p).$$

To learn a less-biased selection probability, we use Gumbel-Softmax:

$$\left(\bar{\Pi}_{i,j}^p\right)_o = \frac{\exp\left(\frac{\log \alpha_o + G_o}{\tau}\right)}{\sum_{o' \in \{0,1\}} \exp\left(\frac{\log \alpha_{o'} + G_{o'}}{\tau}\right)} \text{ where } o \in \{0,1\}.$$

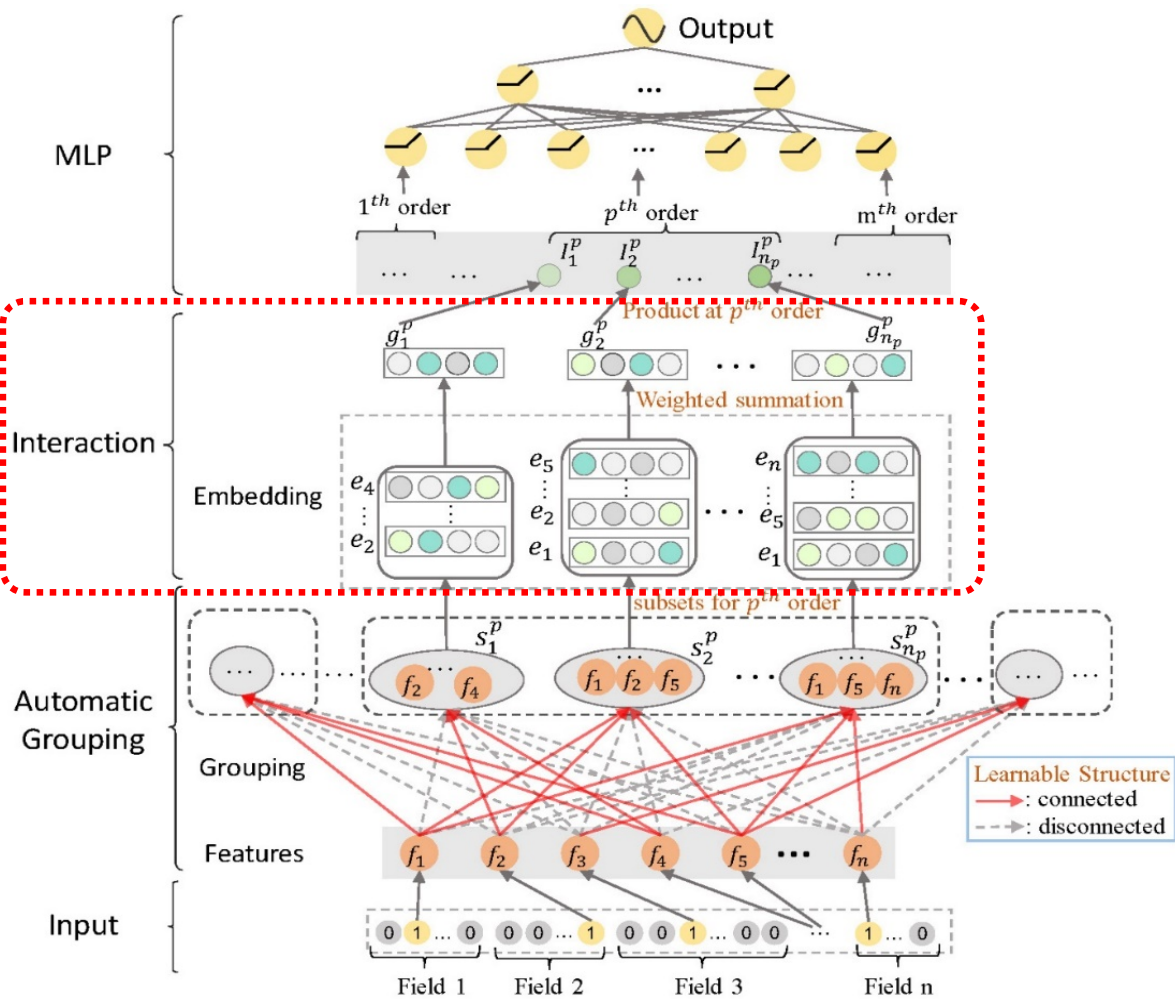
$$\alpha_0 = \frac{1}{1 + \exp(-\alpha_{i,j}^p)} \quad \alpha_1 = \frac{\exp(-\alpha_{i,j}^p)}{1 + \exp(-\alpha_{i,j}^p)}$$

$$G_o = -\log(-\log u) \text{ where } u \sim \text{Uniform}(0,1)$$

**Trainable Parameters:**  $\{\alpha_{i,j}^p\}$



# Interaction Stage



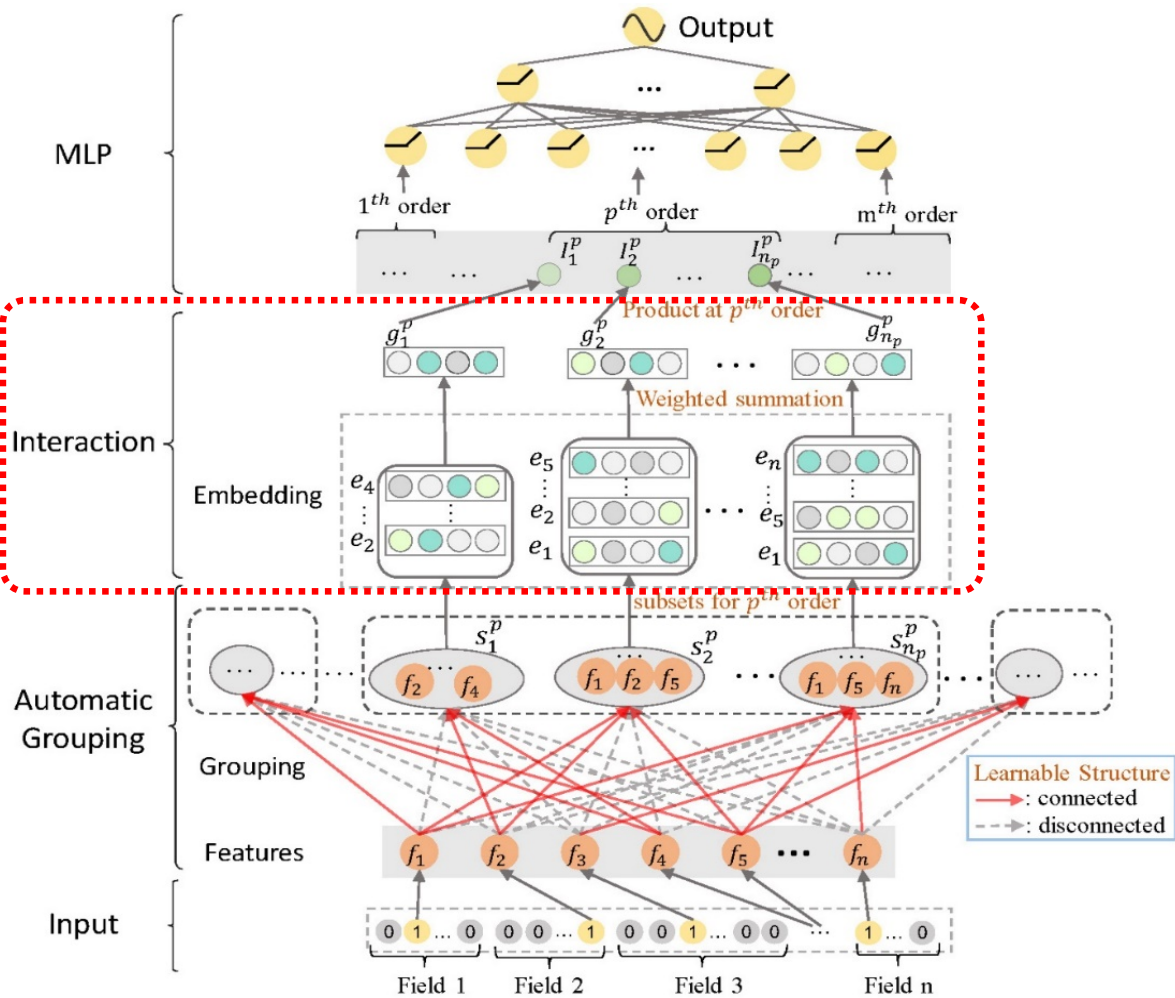
Feature set representation:

$$g_j^p = \sum_{f_i \in S_j^p} w_i^p e_i$$

$S_j^p$ : the  $j^{\text{th}}$  feature set for order- $p$  feature interactions.

$e_i$ : embedding for feature  $f_i$

$w_i^p$ : weights of embeddings in feature set  $S_j^p$ .



Interaction at a given order:

- Inspired by the reformulation of FM:

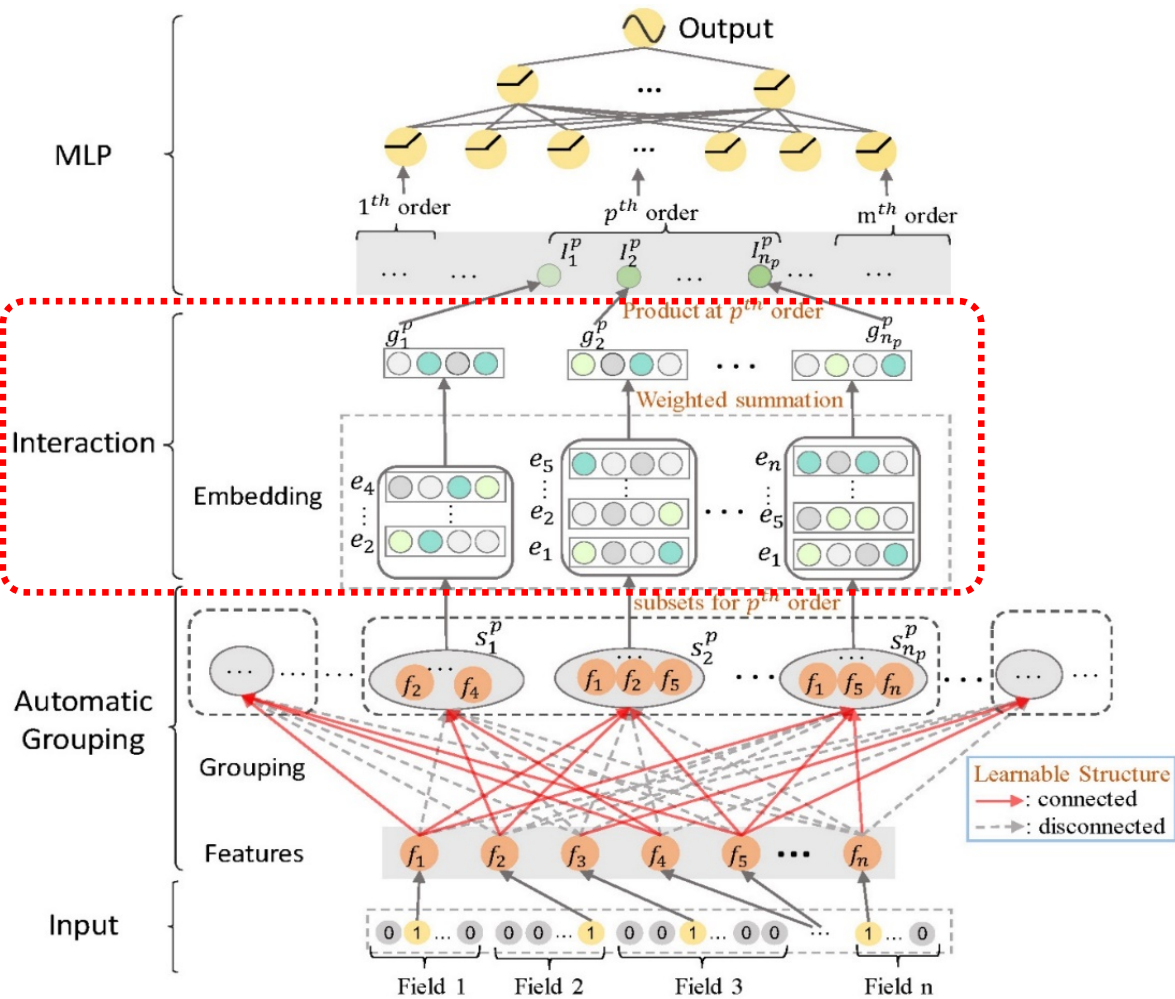
$$\hat{y}(x) = w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle e_i, e_j \rangle x_i x_j$$

$$= w_0 + \sum_{i=1}^n w_i x_i + \frac{1}{2} \left( \left( \sum_{i=1}^n x_i e_i \right)^2 - \sum_{i=1}^n (x_i e_i)^2 \right)$$

- The order- $p$  interaction in a given set  $s_j^p$  is defined as:

$$I_j^p = \begin{cases} \left( g_j^p \right)^p - \sum_{f_i \in s_j^k} (w_i^p e_i)^p \in R, p \geq 2 \\ g_j^p \in R^k, & p = 1 \end{cases}$$

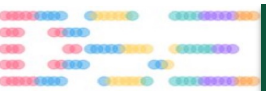
# Interaction Stage



Parameter Training:  
The structural parameters  $\{\alpha_{i,j}^p\}$  and other normal parameters (embedding parameters and network parameters) are optimized alternatively in bi-level optimization (DARTS).



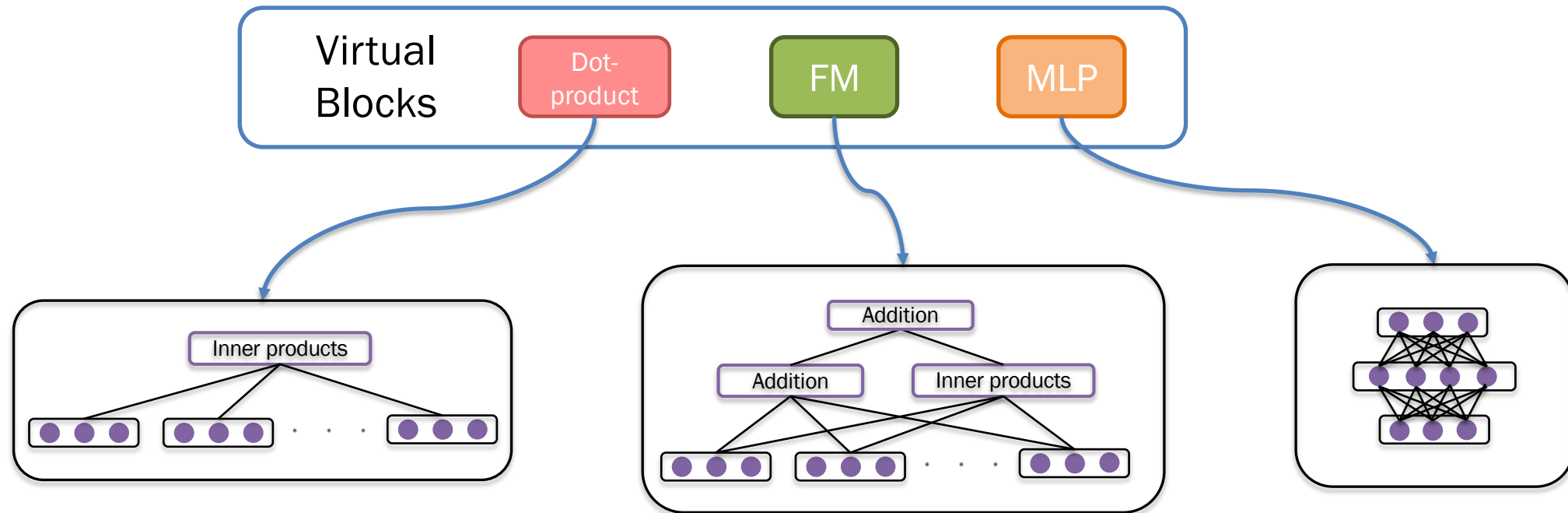
- **AutoML in Embedding Layer**
  - **NIS** - Neural Input Search for Large Scale Recommendation Models (KDD'2020)
  - **ESAPN** - Automated Embedding Size Search in Deep Recommender Systems (SIGIR'2020)
  - **AutoDim** - Field-aware Embedding Dimension Search in Recommender Systems (WWW'2021)
  - **AutoDis** - Automatic Discretization for Embedding Numerical Features in CTR Prediction (AAAI'2021)
- **AutoML in Interaction Layer**
  - **AutoFIS** - Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction (KDD'2020)
  - **AutoGroup** - Automatic Feature Grouping for Modelling Explicit High-Order Feature Interactions in CTR Prediction (SIGIR'2020)
  - **AutoCTR** - Towards Automated Neural Interaction Discovery for Click-Through Rate Prediction (KDD'2020)



# AutoCTR - Hierarchical Search Space

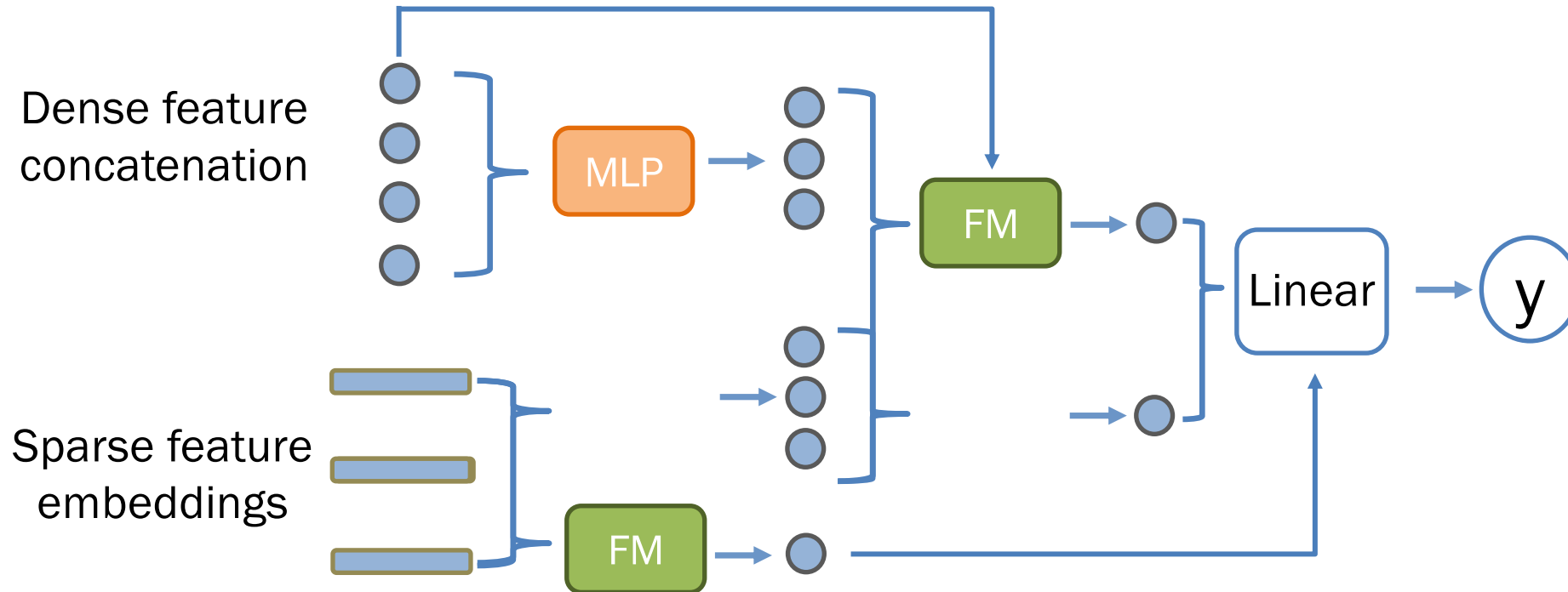


- Virtual block abstraction
  - Properties: functionality complementary, complexity aware, ...
  - Examples: MLP block, dot-product block, factorization-machine block, ...

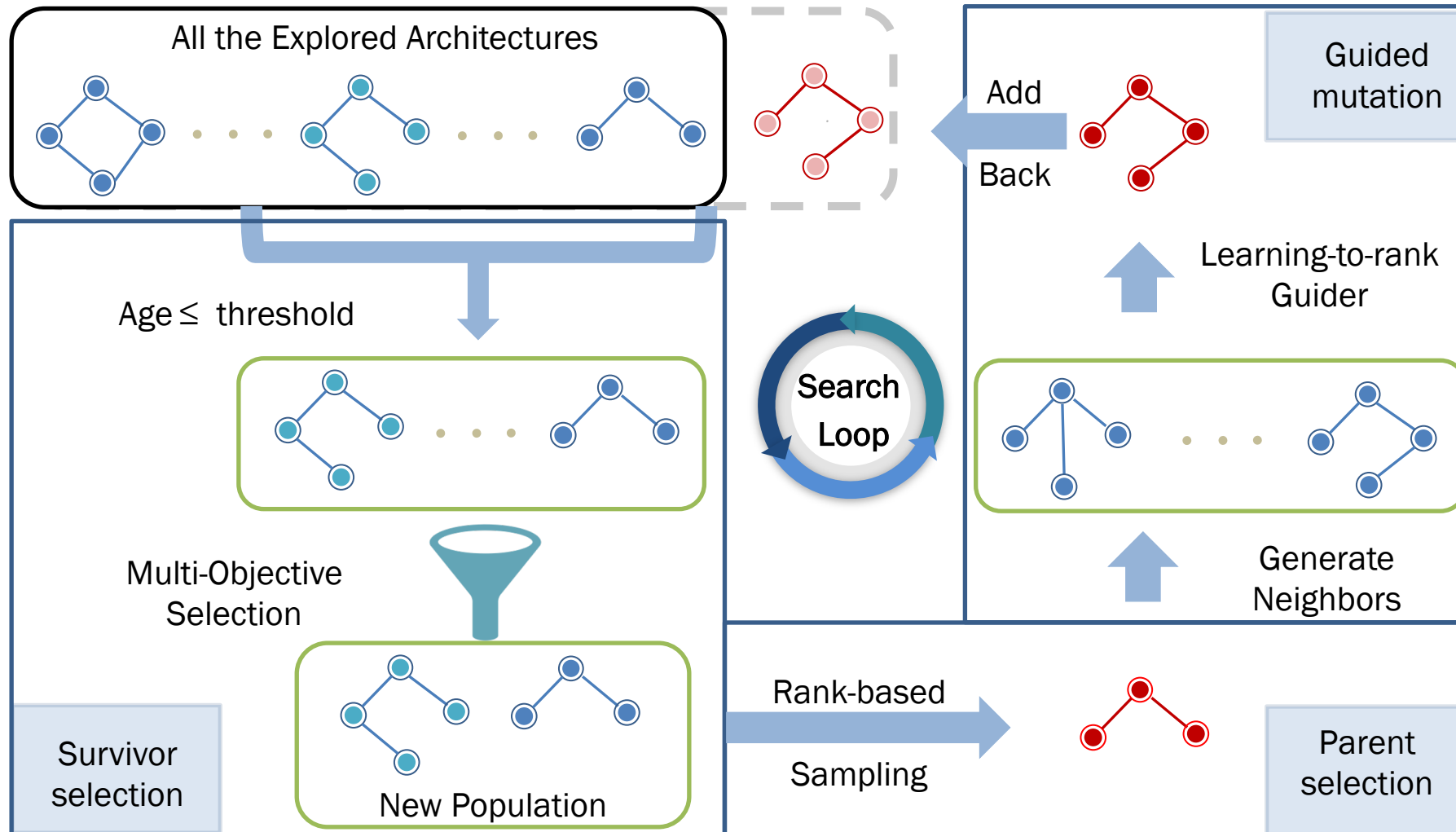


## Search space construction

- DAG of virtual blocks and grouped feature embeddings
- Both block hyperparameters and connection among blocks are to be searched

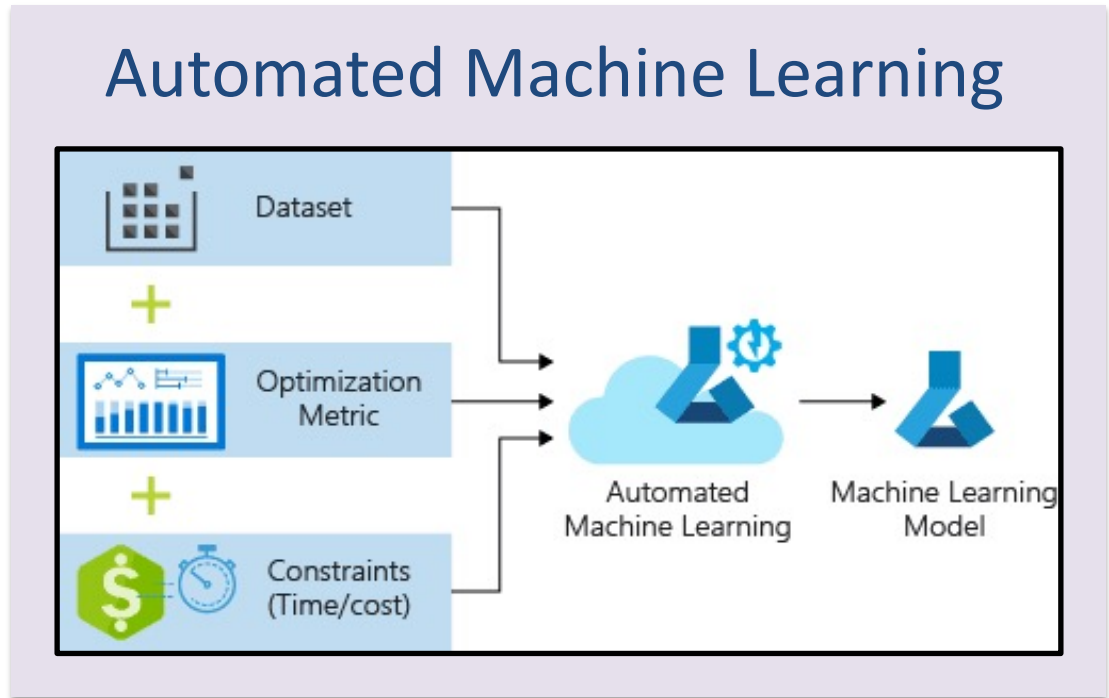


# AutoCTR - Multi-Objective Evolutionary Search Algorithm



# Conclusion

- Deep architectures are designed by the machine automatically
- Advantages
  - Less expert knowledge
  - Saving time and efforts
  - Different data → different architectures



# Future Directions

- Applying AutoML to more tasks
  - Feature engineering, model selection, optimization algorithm, model evaluation, etc

