



# Automated Machine Learning for Recommendations: Fundamentals and Advances

Xiangyu Zhao

City University of Hong Kong

[xianzhao@cityu.edu.hk](mailto:xianzhao@cityu.edu.hk)

Wenqi Fan

The Hong Kong Polytechnic University

[wenqifan03@gmail.com](mailto:wenqifan03@gmail.com)

Huifeng Guo

Huawei Noah's Ark Lab

[huifeng.guo@huawei.com](mailto:huifeng.guo@huawei.com)

Bo Chen

Huawei Noah's Ark Lab

[chenbo116@huawei.com](mailto:chenbo116@huawei.com)

Yejing Wang

City University of Hong Kong

[adave631@gmail.com](mailto:adave631@gmail.com)

Ruiming Tang

Huawei Noah's Ark Lab

[tangruiming@huawei.com](mailto:tangruiming@huawei.com)

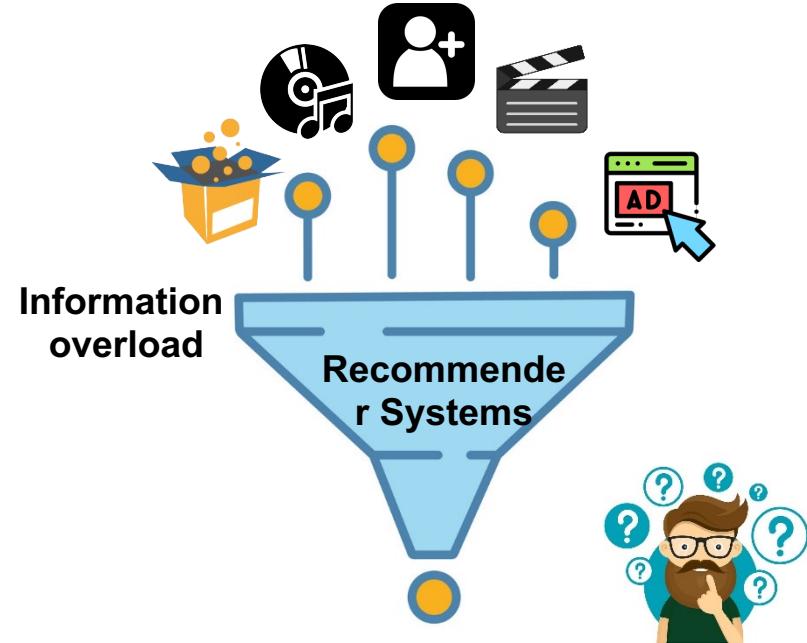


# Table of Contents

- Introduction
  - Background: Deep Recommender Systems
- Preliminary of AutoML
- DRS Embedding Components
  - Single Embedding Search
  - Group Embedding Search
- DRS Interaction Components
  - Feature Interaction Search
  - Interaction Function Search
  - Interaction Block Search
- DRS Comprehensive Search & System
- Conclusion & Future Direction

# Recommender Systems

## Age of Information



Recommend item X to

**Items** can be: Products, News, Movies, Videos, Friends, etc.

# Recommender Systems



Recommendation has been widely applied in online services:  
- E-commerce, Content Sharing, Social Networking ...

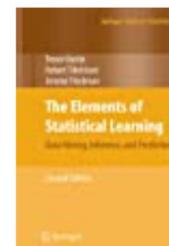


## Product Recommendation

Frequently bought together



+



+



A

B

C

Total price: \$208.9

Add all three to Cart

Add all three to List

# Recommender Systems



Recommendation has been widely applied in online services:  
- E-commerce, **Content Sharing**, Social Networking ...

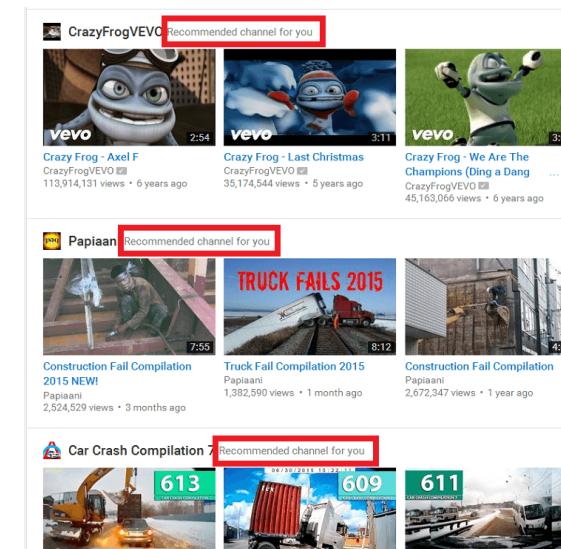


## News/Video/Image Recommendation

For you  
Recommended based on your interests

This Research Paper From Google Research Proposes A 'Message Passing Graph Neural Network' That Explicitly Models Spatio-Temporal Relations  
MarkTechPost · 2 days ago

Tested: Brydge MacBook Vertical Dock, completing my MacBook Pro desktop  
9to5Mac · 21 hours ago



# Recommender Systems



Recommendation has been widely applied in online services:  
- E-commerce, Content Sharing, **Social Networking** ...

facebook

LinkedIn®



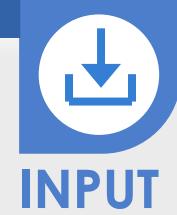
## Friend Recommendation

The screenshot shows a Facebook user profile for Andrew Torba. On the left, there's a sidebar with 'FAVORITES' including 'News Feed' (selected), 'Messages', 'Events' (2), 'Find Friends' (17), 'Tech.li', 'Kuhcoon', and 'APX'. Below that is a 'PAGES' section with several blacked-out links. A modal window titled 'Are They Your Friends Too?' appears, stating 'These people now have 1 or more friends in common with you.' It lists four profiles with their mutual friend counts and 'Add Friend' buttons:

Profile Picture	Mutual Friends	Action
	1 mutual friend	Add Friend
	67 mutual friends	Add Friend
	39 mutual friends	Add Friend
	47 mutual friends	Add Friend

At the bottom of the modal is a 'See All Suggestions' button.

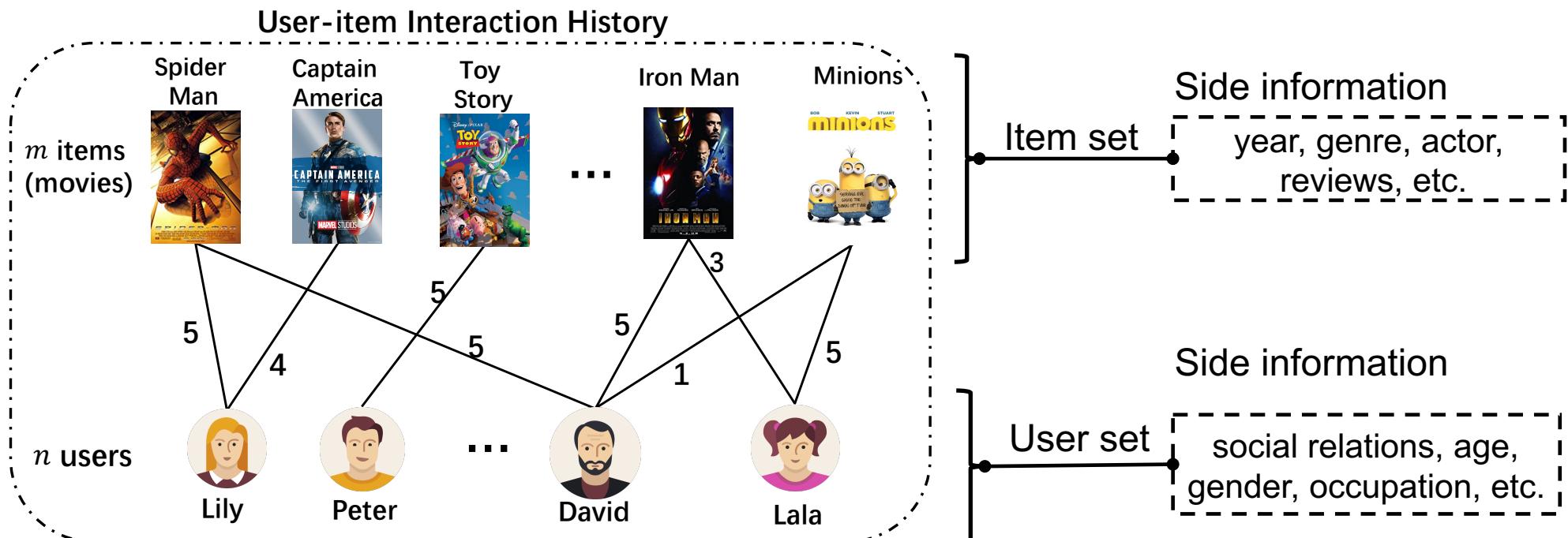
# Problem Formulation



Historical user-item interactions or additional side information (e.g., social relations, item's knowledge, etc.)

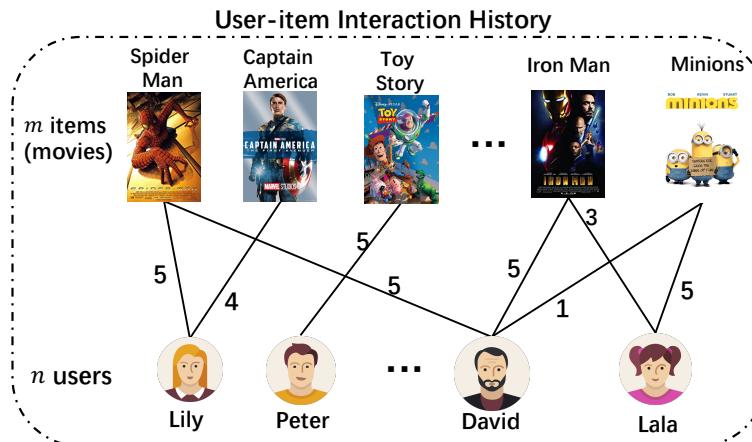


Predict how likely a user would interact with a target Item (e.g., click, view, or purchase)



# Recommender Systems

- Collaborative Filtering (CF) is the most well-known technique for recommendation.
  - Similar users (with respect to their historical interactions) have similar preferences.
  - Modelling users' preference on items based on their past interactions (e.g., ratings and clicks).
- Learning representations of users and items is the key of CF.

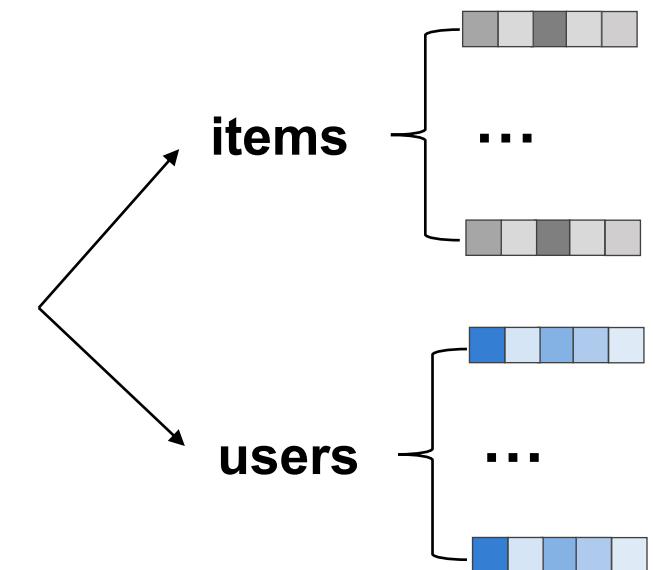


User-item Rating Matrix

	Spider Man	Captain America	Toy Story	Iron Man	Minions	R
Lily	5	4	?		?	?
Peter	?	?	5	...	?	?
David				...		
Lala				1		5
	5	?	?		5	1
	?	?	?	...	2	5

m items (movies)

n users

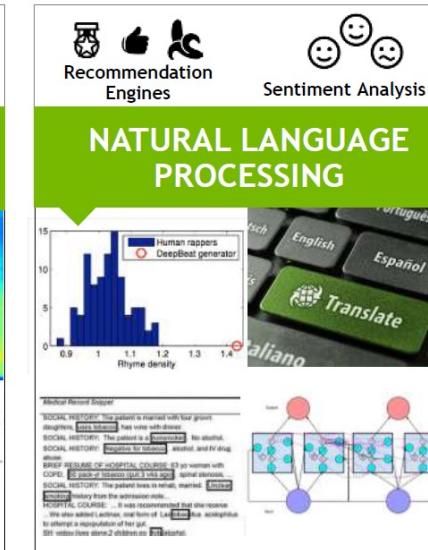
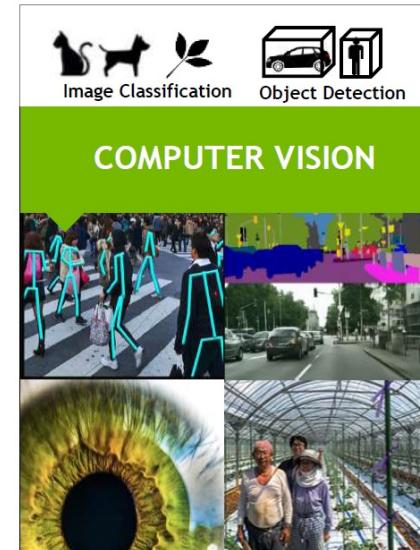
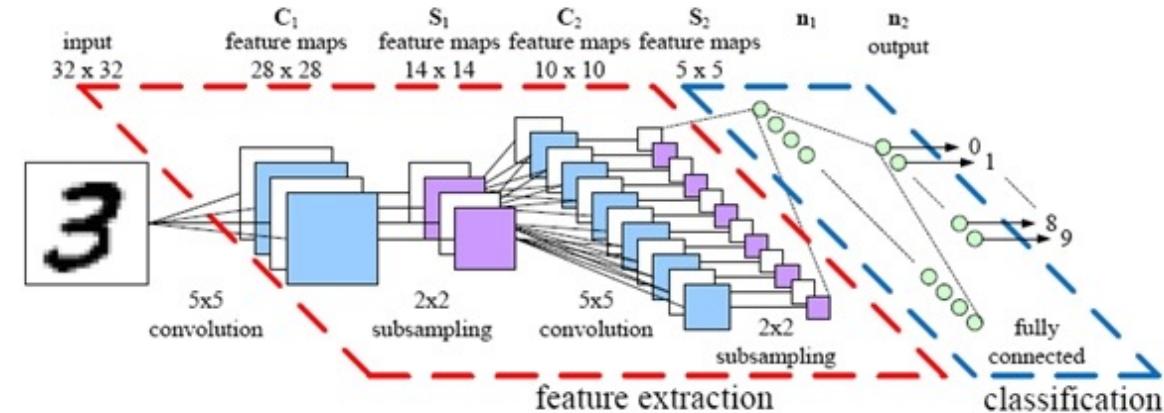
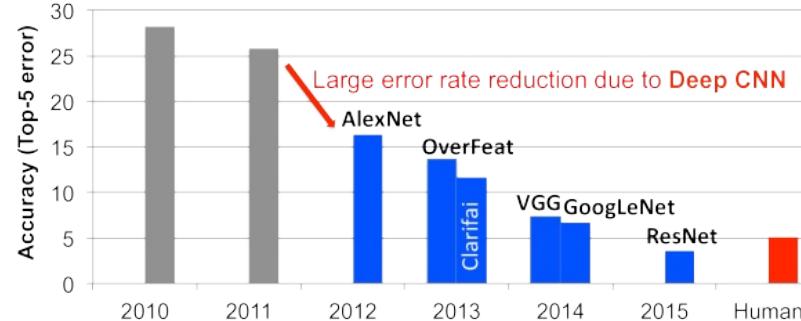


Task: predicting missing movie ratings in Netflix.

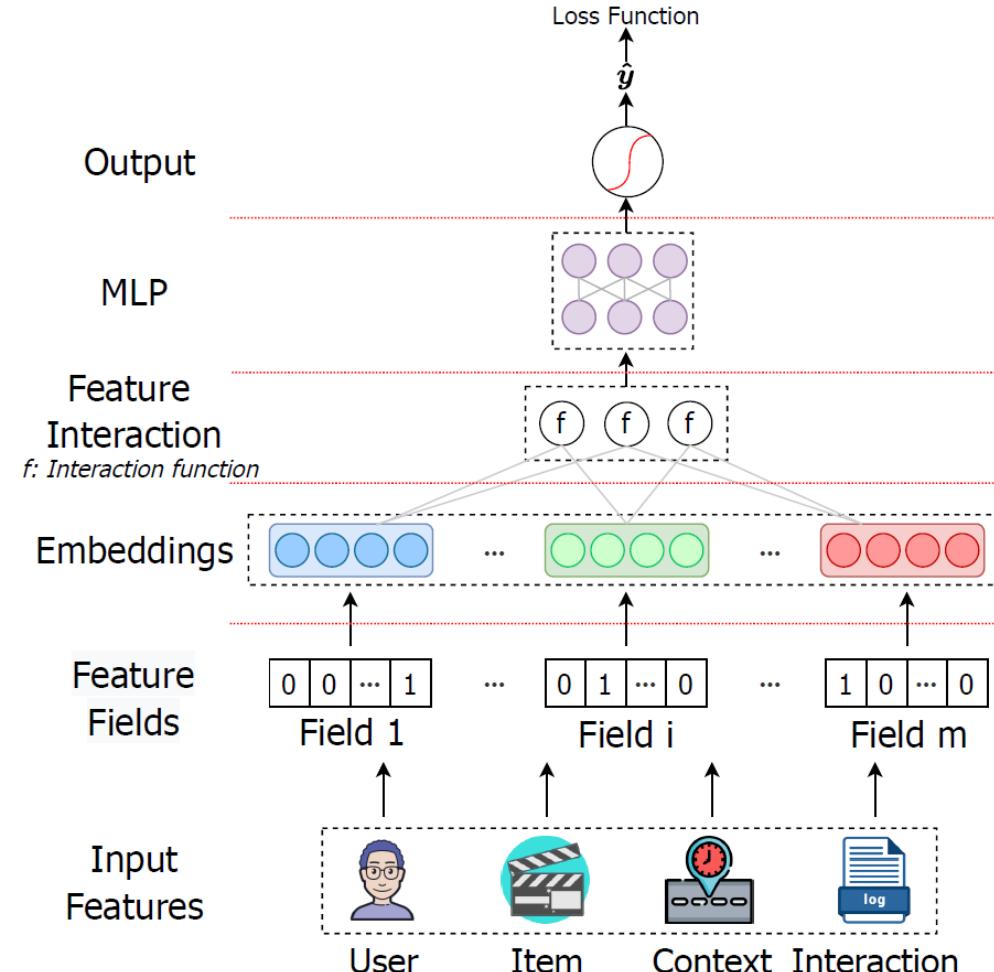
# Deep Learning is Changing Our Lives



## IMAGENET



# Recommender Systems Architecture

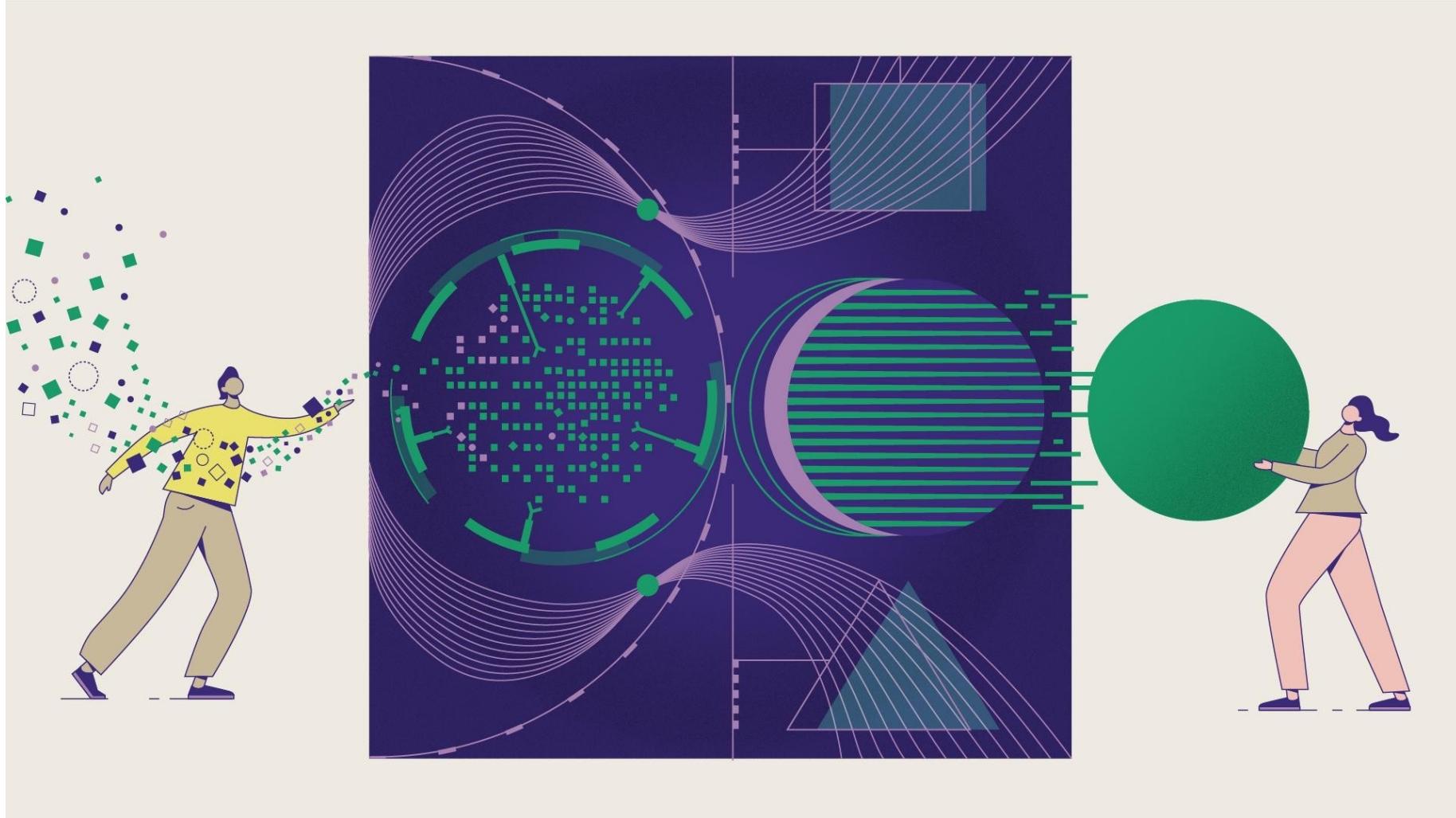




# Table of Contents

- Introduction
  - Background: Deep Recommender Systems
- **Preliminary of AutoML**
- DRS Embedding Components
  - Single Embedding Search
  - Group Embedding Search
- DRS Interaction Components
  - Feature Interaction Search
  - Interaction Function Search
  - Interaction Block Search
- DRS Comprehensive Search & System
- Conclusion & Future Direction

# Why AutoML?



# Success of Machine Learning



Astronomy

Robotic

Creative  
Arts

Teaching

Material  
Design

Energy

Game Play

Search

Chemistry

Image  
Recognition

Weather  
Prediction

Health  
Care

Physics

Manufacturing

Service

Product  
Recommendation

Drug  
Discovery

Maintenance  
Prediction

Traffic  
Prediction

Retail

Financial  
Services

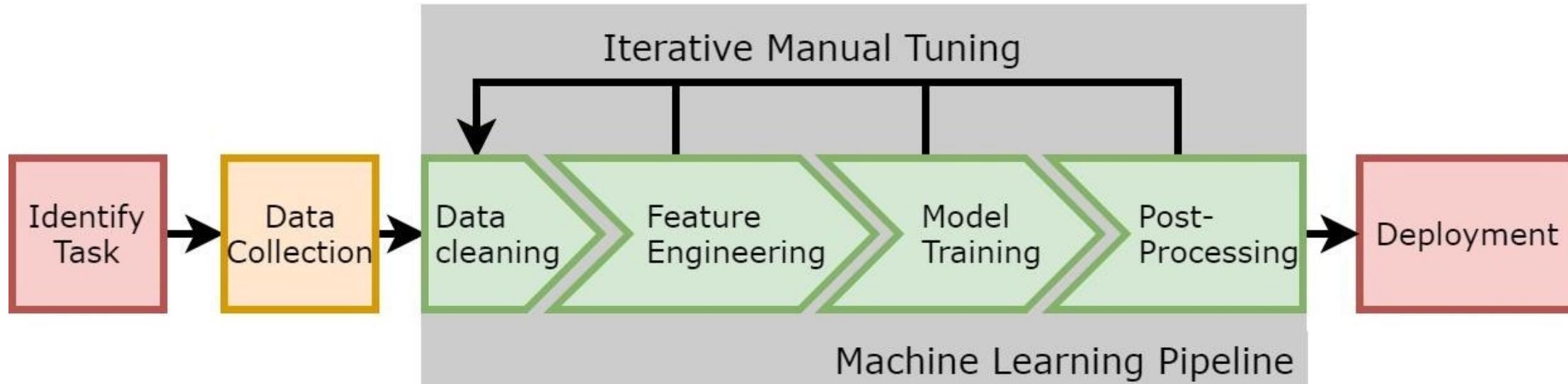
Credit  
Assignment

Social  
Media

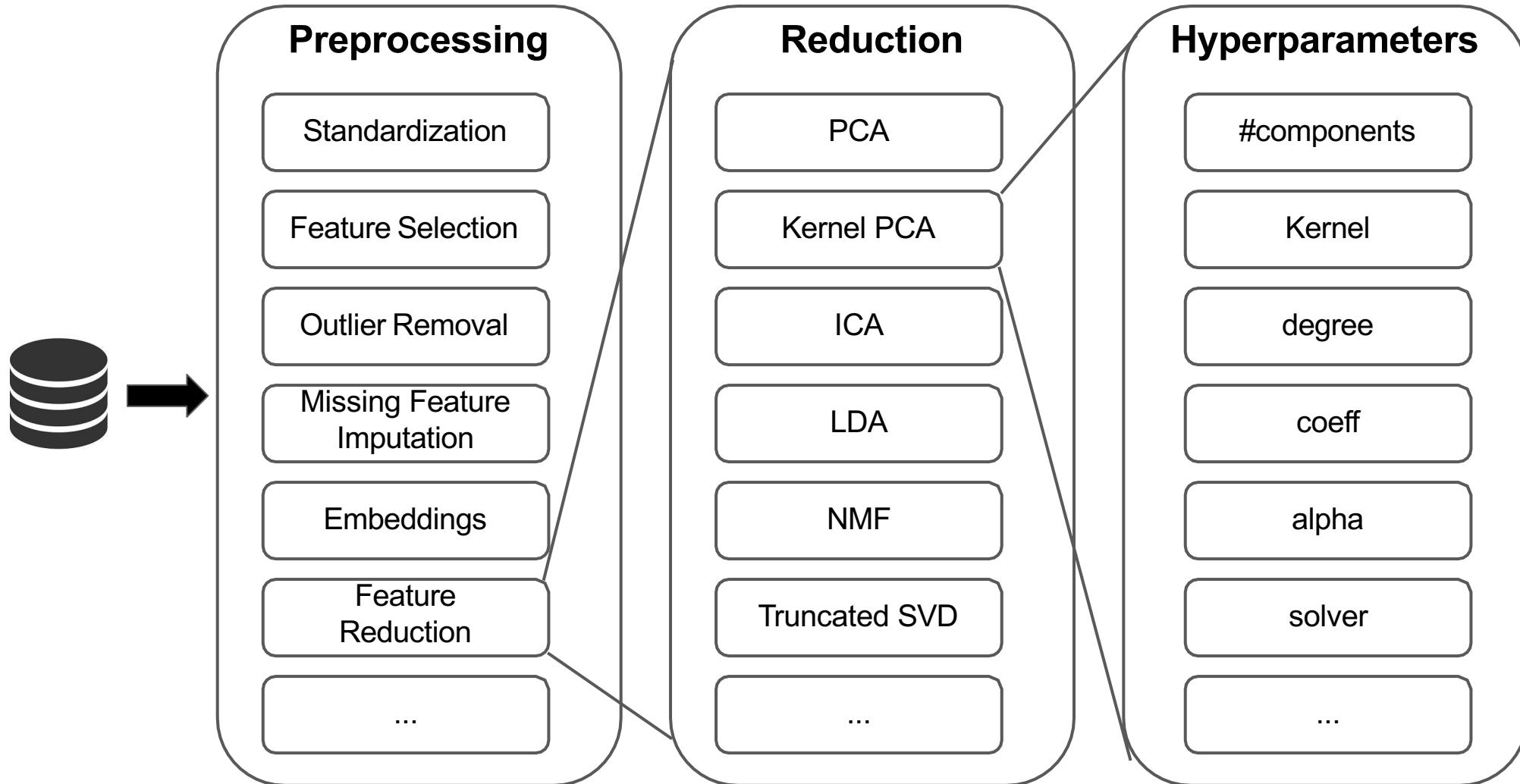
Media

Summary  
Generation

# Machine Learning Pipeline

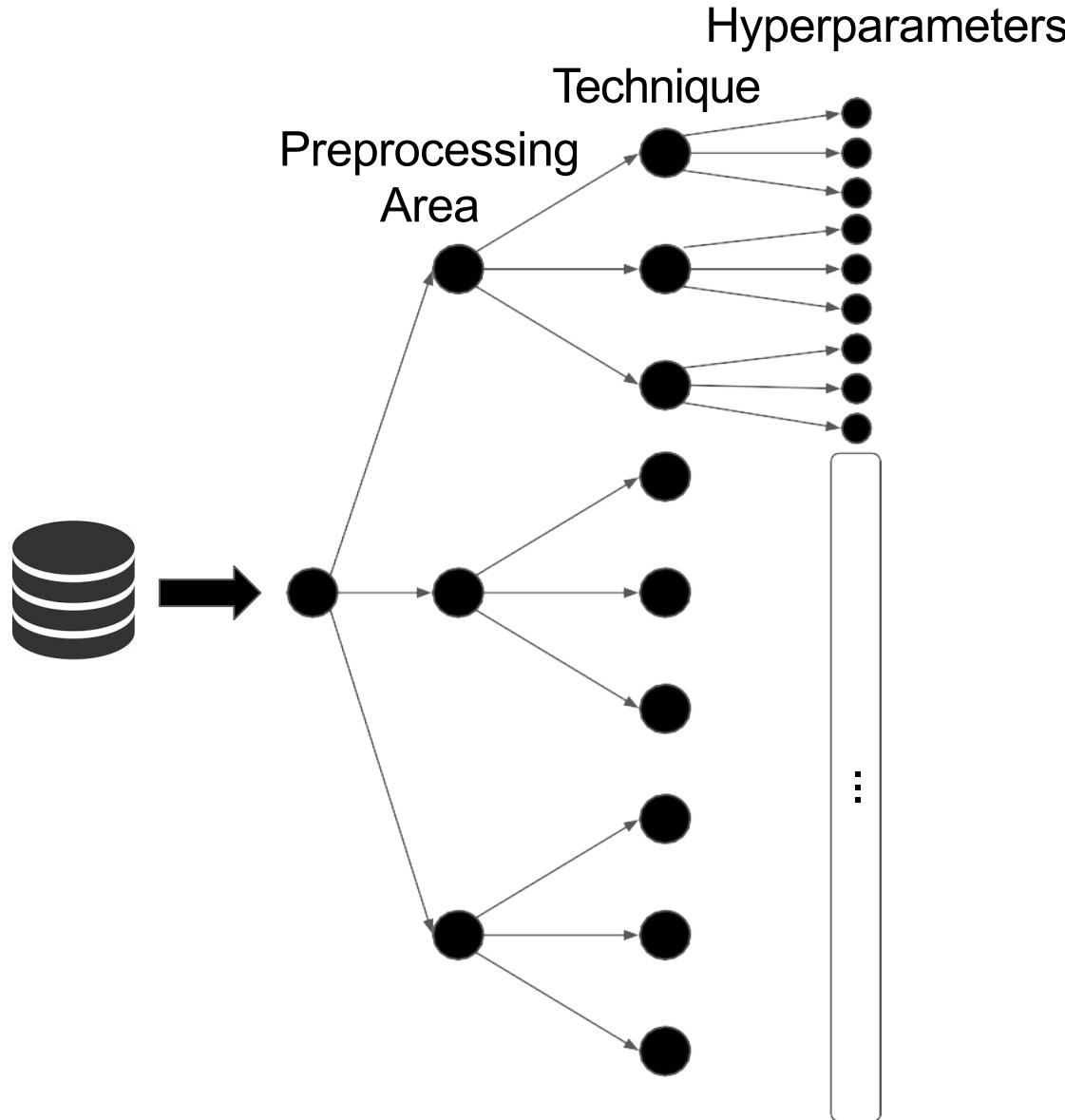


# Preprocessing?



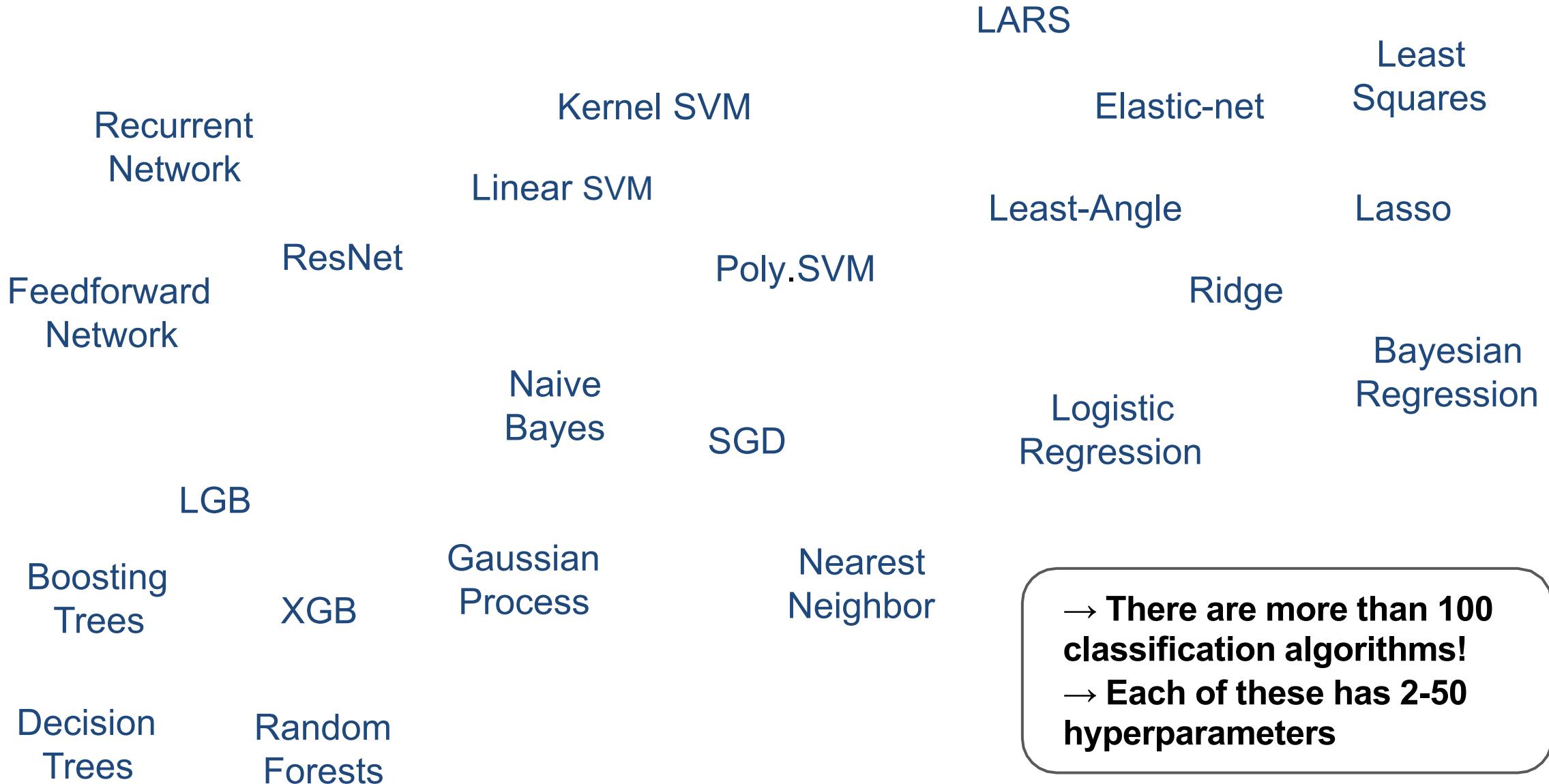
→ We might want more than 1 data preprocessor!

# Complexity of the Preprocessing



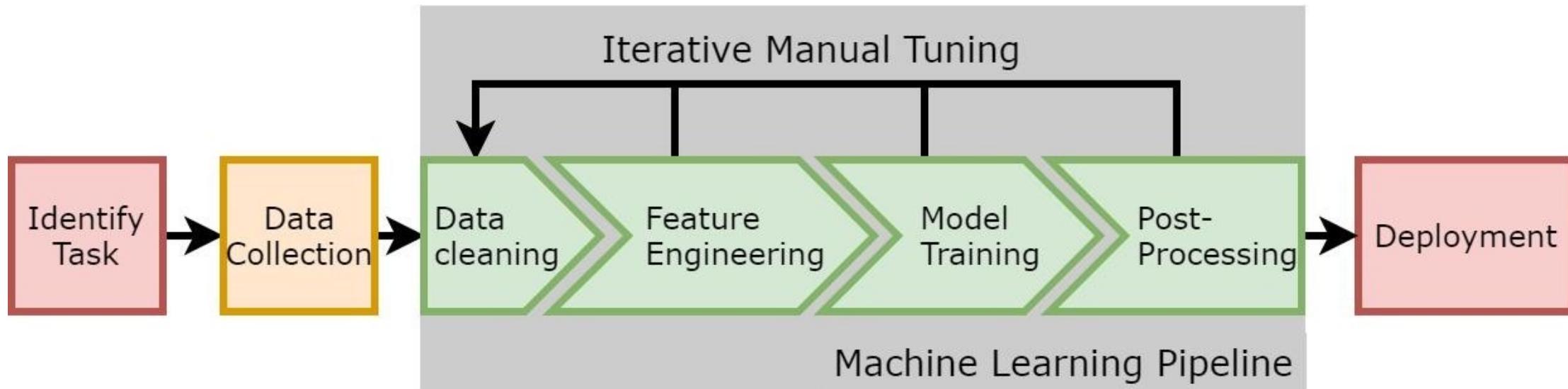
- Naive Assumptions:  
only 3 decisions at each level
- **Possible options:**  $3 \times 3 \times 3 = 27$
- More realistic assumption:  
at least 10 decisions at each level
- **Possible options:**  $10 \times 10 \times 10 = 1000$
- Choose 3 preprocessors instead of 1  
 $\rightarrow 1000 \times 1000 \times 1000 =$   
**1 000 000 000**
- Still naive!  
 $\rightarrow$  Hyperparameters are often continuous and not discrete  
 $\rightarrow$  **infinite amount of settings!**

# Classification Algorithms

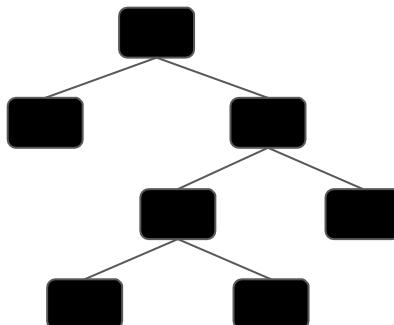


→ There are more than 100 classification algorithms!  
→ Each of these has 2-50 hyperparameters

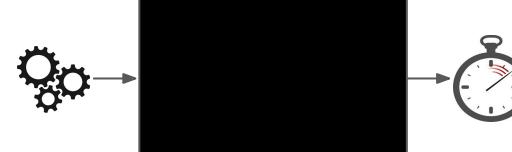
# Challenges in Designing ML Pipelines



Complex  
Search Space



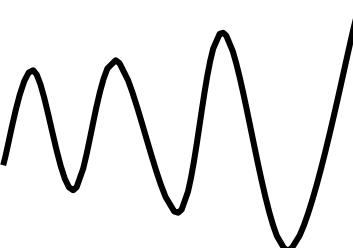
Black-Box  
Problem



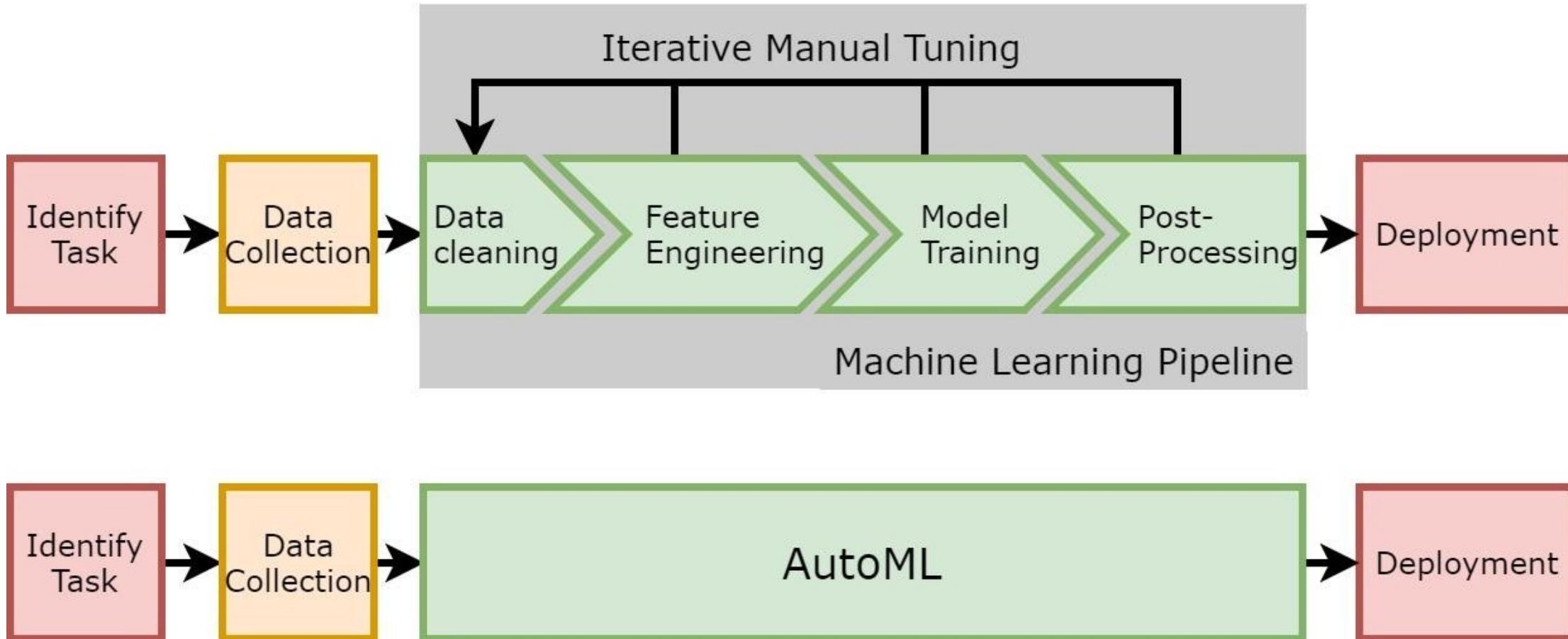
Expensive  
Evaluations



Noise on  
observations



# From Manual ML to Automated ML

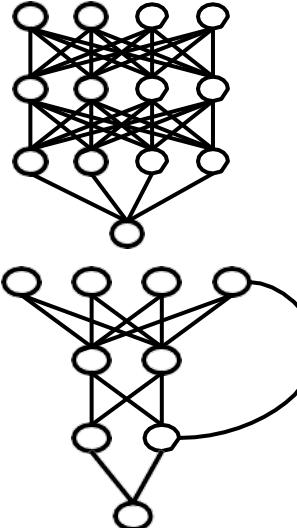


# Design Decisions by AutoML



classifier	#λ
AdaBoost (AB)	4
Bernoulli naïve Bayes	2
decision tree (DT)	4
extremal. rand. trees	5
Gaussian naïve Bayes	-
gradient boosting (GB)	6
kNN	3
LDA	4
linear SVM	4
kernel SVM	7
multinomial naïve Bayes	2
passive aggressive	3
QDA	2
random forest (RF)	5
Linear Class. (SGD)	10

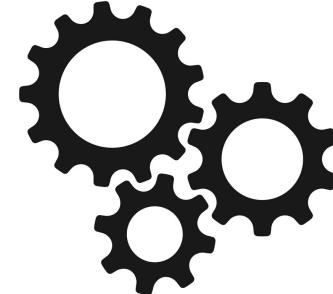
Algorithms



Architecture Design

preprocessor	#λ
extremal. rand. trees prepr.	5
fast ICA	4
feature agglomeration	4
kernel PCA	5
rand. kitchen sinks	2
linear SVM prepr.	3
no preprocessing	-
nystroem sampler	5
PCA	2
polynomial	3
random trees embed.	4
select percentile	2
select rates	3
one-hot encoding	2
imputation	1
balancing	1
rescaling	1

Pre-processing



Hyper-parameters

...

# Neural Architecture Search (NAS)



- Find neural architecture  $A$  such that deep learning works best for given data
  - Measured by validation error of architecture  $A$  with trained weights  $w^*(A)$

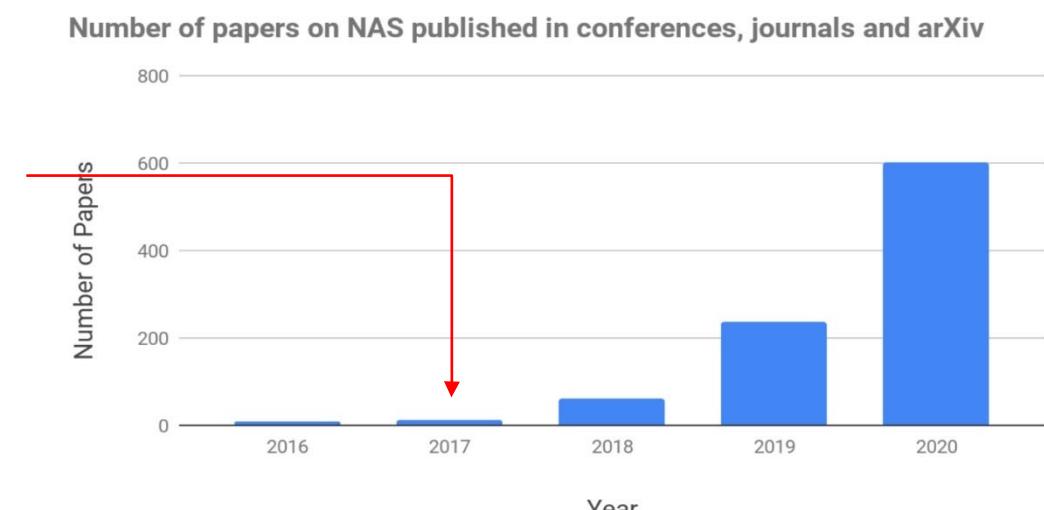
$$\min_{A \in \mathcal{A}} \mathcal{L}_{\text{val}}(w^*(A), A)$$

$$\text{s.t. } w^*(A) \in \operatorname{argmin}_w \mathcal{L}_{\text{train}}(w, A)$$

- Famously tackled by

reinforcement learning [[Zoph & Le, ICLR 2017](#)]

- 12.800 architectures trained fully
- 800 GPUs for 2 weeks (about \$60.000 USD)



# NAS is Exploding!



Journal of Machine Learning Research 20 (2019) 1-21

Submitted 9/18; Revised 3/19; Published 3/19

## Neural Architecture Search: A Survey

**Thomas Elsken**

*Bosch Center for Artificial Intelligence  
71272 Renningen, Germany  
and University of Freiburg*

THOMAS.ELSKEN@DE.BOSCH.COM

**Jan Hendrik Metzen**

*Bosch Center for Artificial Intelligence  
71272 Renningen, Germany*

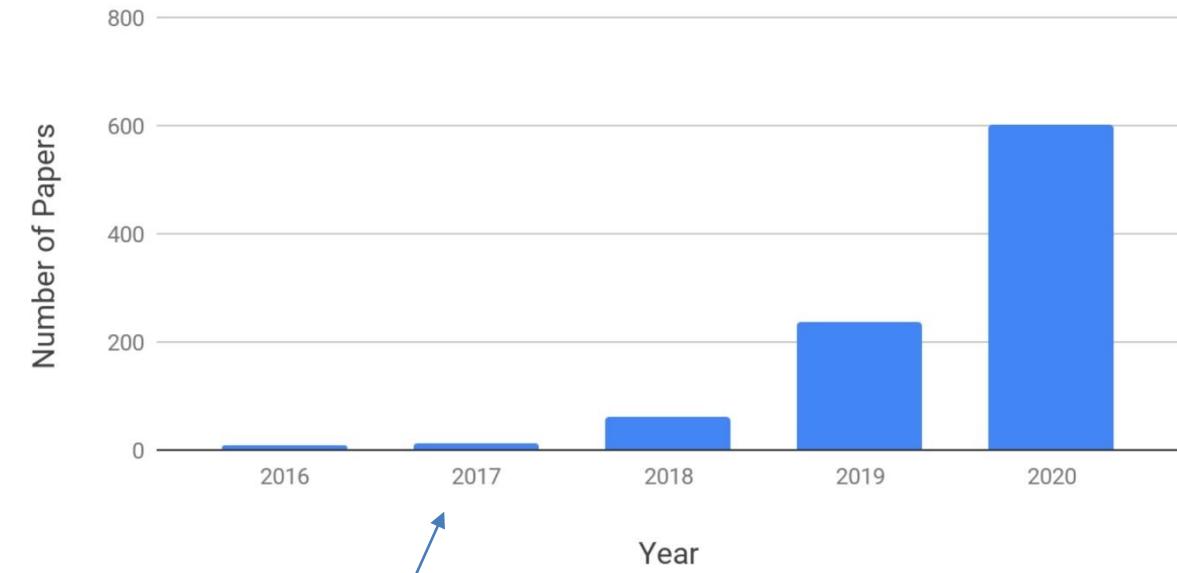
JANHENDRIK.METZEN@DE.BOSCH.COM

**Frank Hutter**

*University of Freiburg  
79110 Freiburg, Germany*

FH@CS.UNI-FREIBURG.DE

Number of papers on NAS published in conferences, journals and arXiv



[Neural Architecture Search with Reinforcement Learning](#)  
[\[Zoph & Le, ICLR 2017\]](#)

# Framework: Trial and Update

- Almost all NAS algorithms are based on the “trial and update” framework
  - Starting with a set of initial architectures (e.g., manually defined) as individuals
  - Assuming that better architectures can be obtained by slight modification
  - Applying different operations on the existing architectures
  - Preserving the high-quality individuals and updating the individual pool
  - Iterating till the end
- Three fundamental requirements
  - The building blocks: defining the search space (dimensionality, complexity, etc.)
  - The representation: defining the transition between individuals
  - The evaluation method: determining if a generated individual is of high quality

# Framework: Building Blocks



- Building blocks are like basic *genes* for these individuals
- Some examples here
  - Genetic CNN: only  $3 \times 3$  convolution is allowed to be searched (followed by default BN and ReLU operations),  $3 \times 3$  pooling is fixed
  - NASNet: 13 operations shown below
  - PNASNet: 8 operations, removing those never-used ones from NASNet
  - ENASNet: 6 operations
  - DARTS: 8 operations
- identity
- $1 \times 7$  then  $7 \times 1$  convolution
- $3 \times 3$  average pooling
- $5 \times 5$  max pooling
- $1 \times 1$  convolution
- $3 \times 3$  depthwise-separable conv
- $7 \times 7$  depthwise-separable conv
- $1 \times 3$  then  $3 \times 1$  convolution
- $3 \times 3$  dilated convolution
- $3 \times 3$  max pooling
- $7 \times 7$  max pooling
- $3 \times 3$  convolution
- $5 \times 5$  depthwise-separable conv

# Framework: Search

- Finding new individuals that have potentials to work better
  - Heuristic search in the large space
- Two mainly applied methods: the **genetic algorithm** and **reinforcement learning**
  - Both are heuristic algorithms applied to the scenarios of a large search space and limited ability to explore every single element in the space
  - A fundamental assumption: both of these heuristic algorithms can preserve good genes and based on which discover possible improvements
- Also, it is possible to integrate architecture search to network optimization
  - These algorithms are often much faster

[Real, 2017] E. Real *et al.*, Large-Scale Evolution of Image Classifiers, *ICML*, 2017.

[Xie, 2017] L. Xie *et al.*, Genetic CNN, *ICCV*, 2017.

[Zoph, 2018] B. Zoph *et al.*, Learning Transferable Architectures for Scalable Image Recognition, *CVPR*, 2018.

[Liu, 2018] C. Liu *et al.*, Progressive Neural Architecture Search, *ECCV*, 2018.

[Pham, 2018] H. Pham *et al.*, Efficient Neural Architecture Search via Parameter Sharing, *ICML*, 2018.

[Liu, 2019] H. Liu *et al.*, DARTS: Differentiable Architecture Search, *ICLR*, 2019.

# Framework: Evaluation



- Evaluation aims at determining which individuals are good and to be preserved
- Conventionally, this was often done by training a network from scratch
  - This is extremely time-consuming, so researchers often train NAS on a small dataset like CIFAR and then transfer the found architecture to larger datasets like ImageNet
  - Even in this way, the training process is really slow: Genetic-CNN requires 17 GPU-days for a single training process, and NAS-RL requires more than 20,000 GPU-days
- Efficient methods were proposed later
  - Ideas include parameter sharing (without the need of re-training everything for each new individual) and using a differentiable architecture (joint optimization)
  - Now, an efficient search process on CIFAR can be reduced to a few GPU-hours, though training the searched architecture on ImageNet is still time-consuming

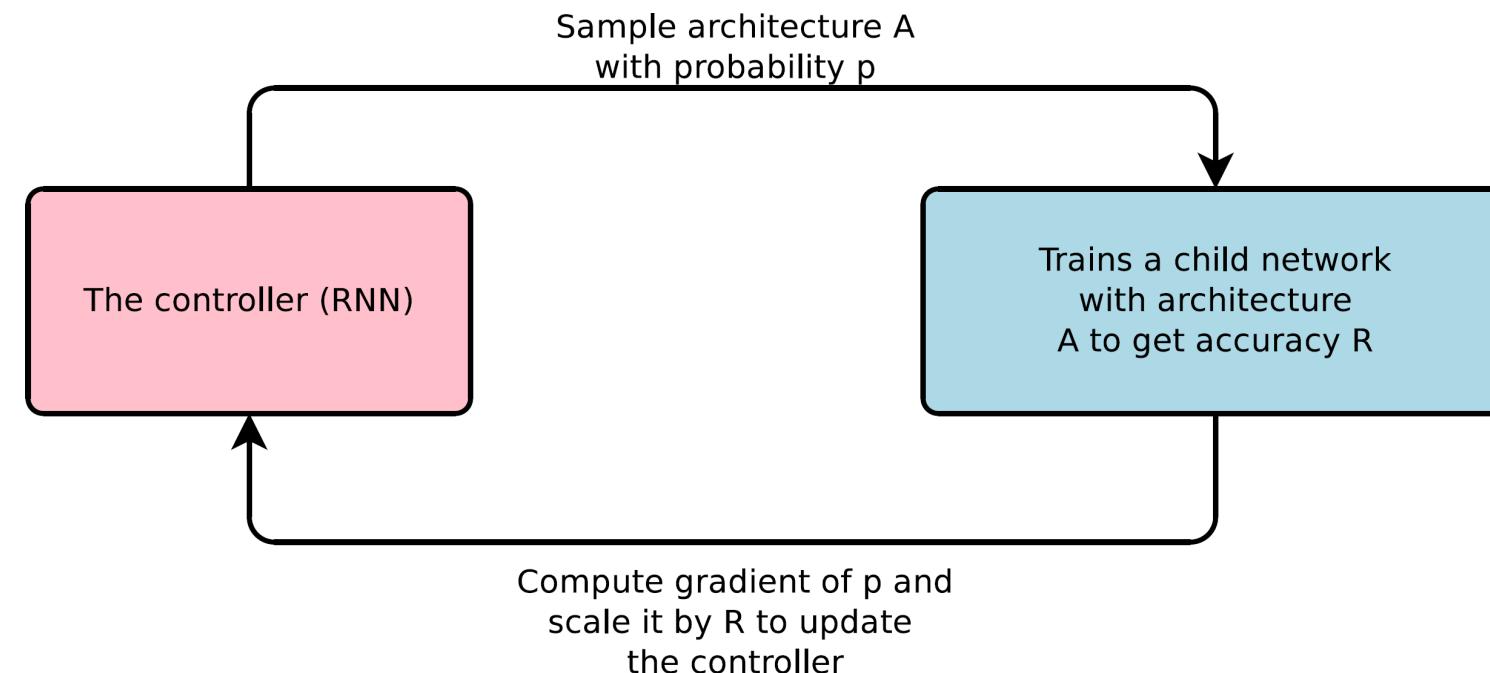
# NAS with Reinforcement Learning



- NAS with Reinforcement Learning [Zoph & Le, ICLR 2017]

- State-of-the-art results for CIFAR-10, Penn Treebank
- Large computational demands:

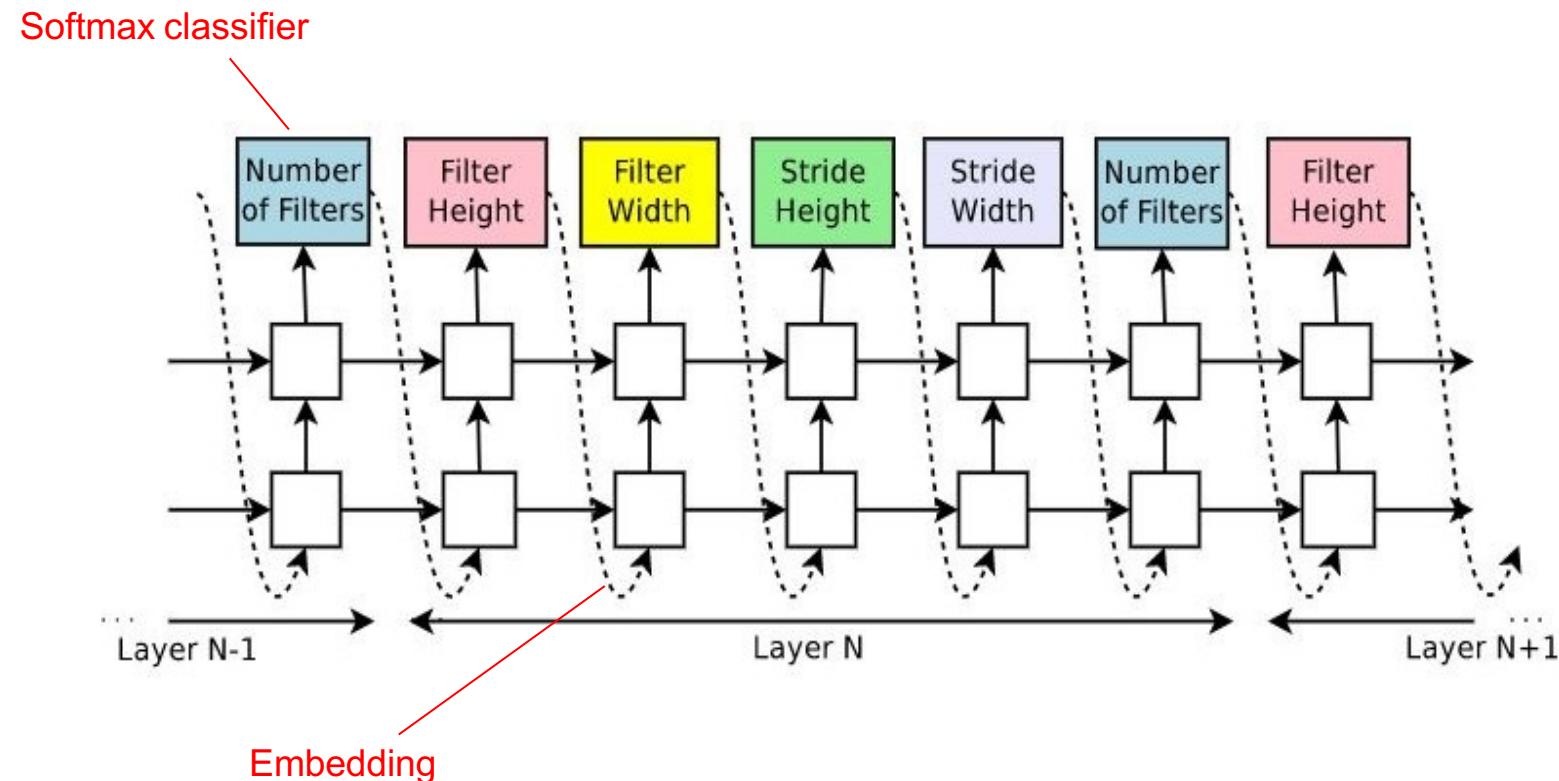
**800 GPUs for 3-4 weeks, 12.800 architectures trained**



# NAS with Reinforcement Learning



- Architecture of neural network represented as string e.g., ["filter height: 5", "filter width: 3", "# of filters: 24"]  
[Zoph & Le, ICLR 2017]
- Controller (RNN) generates string that represents architecture



# Training with REINFORCE



Parameters of Controller RNN

Accuracy of architecture on  
held-out dataset

$$J(\theta_c) = E_{P(a_{1:T}; \theta_c)}[R]$$

Architecture predicted by the controller RNN  
viewed as a sequence of actions

# NAS as Hyperparameter Optimization



- Architecture of neural network represented as string e.g., ["filter height: 5", "filter width: 3", "# of filters: 24"] [Zoph & Le, ICLR 2017]
- We can simply treat these as categorical parameters
  - E.g., 25 cat. parameters for each of the 2 cells in [Zoph et al, CVPR 2018]



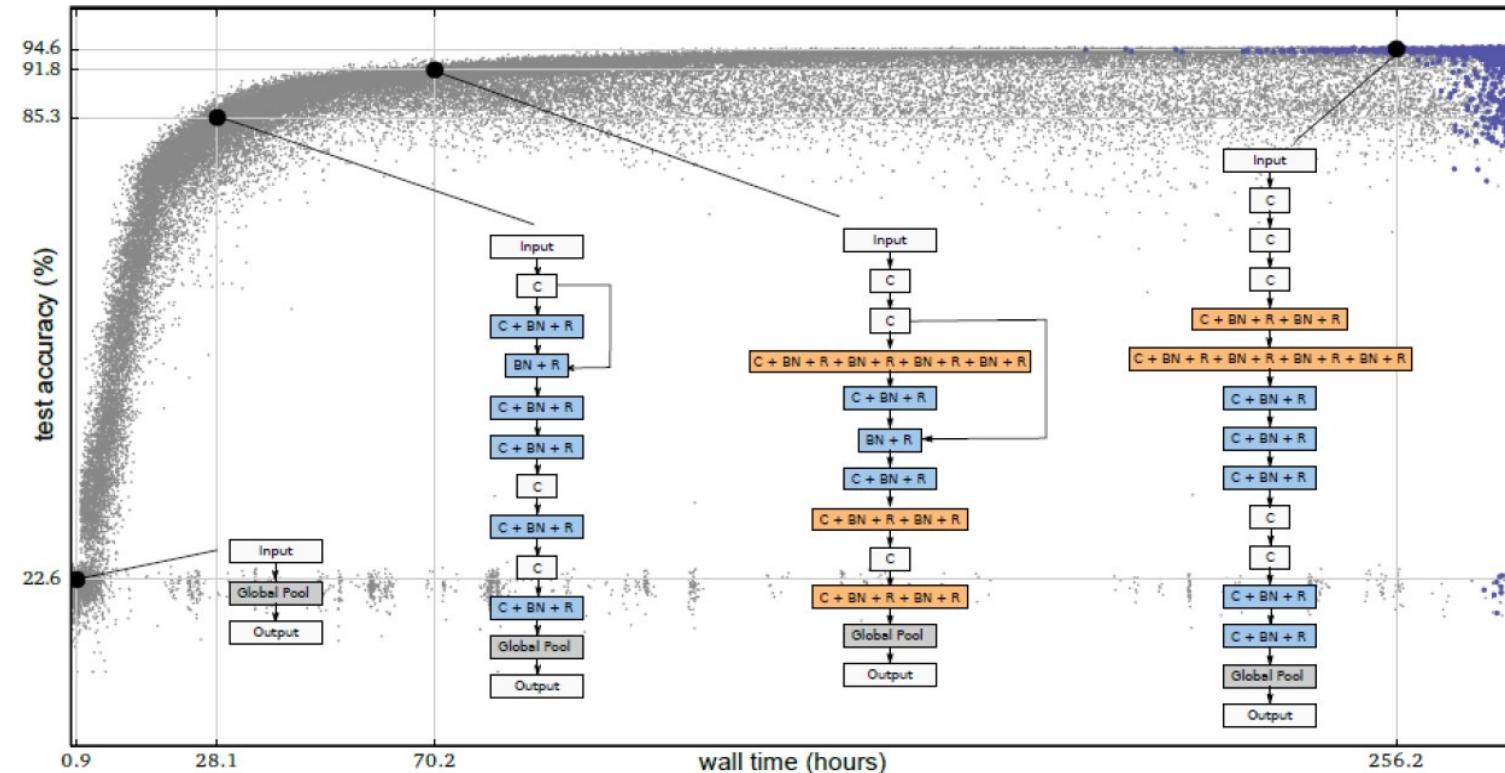
# NAS with Evolution

- Neuroevolution

(already since the 1990s [Angeline et al., 1994; Stanley and Miikkulainen, 2002])

- Mutation steps, such as adding, changing or removing a layer

[Real et al., ICML 2017; Miikkulainen et al., arXiv 2017]

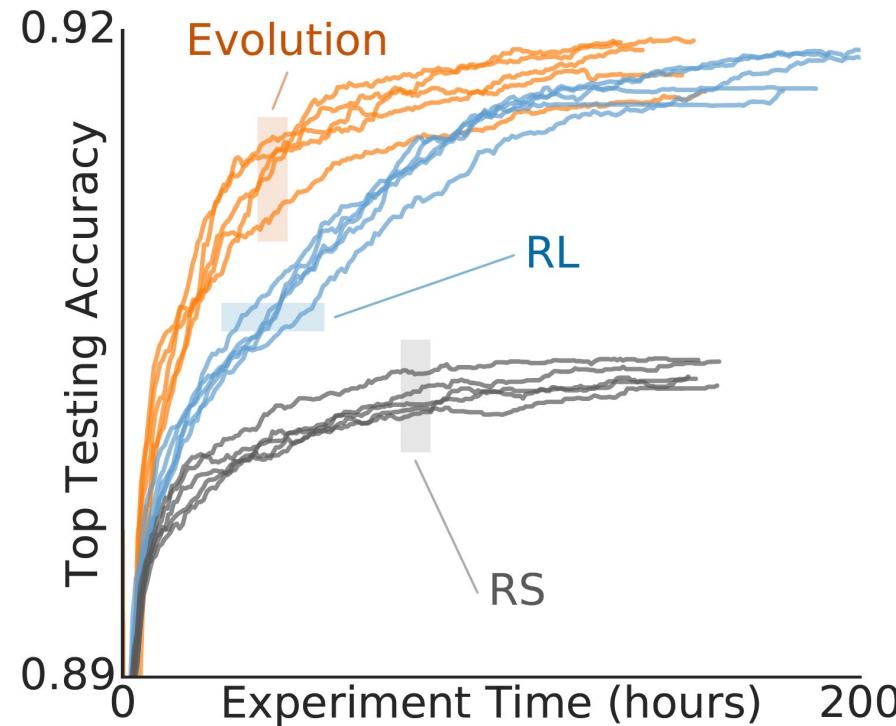


# RL vs. Evolution vs. Random Search

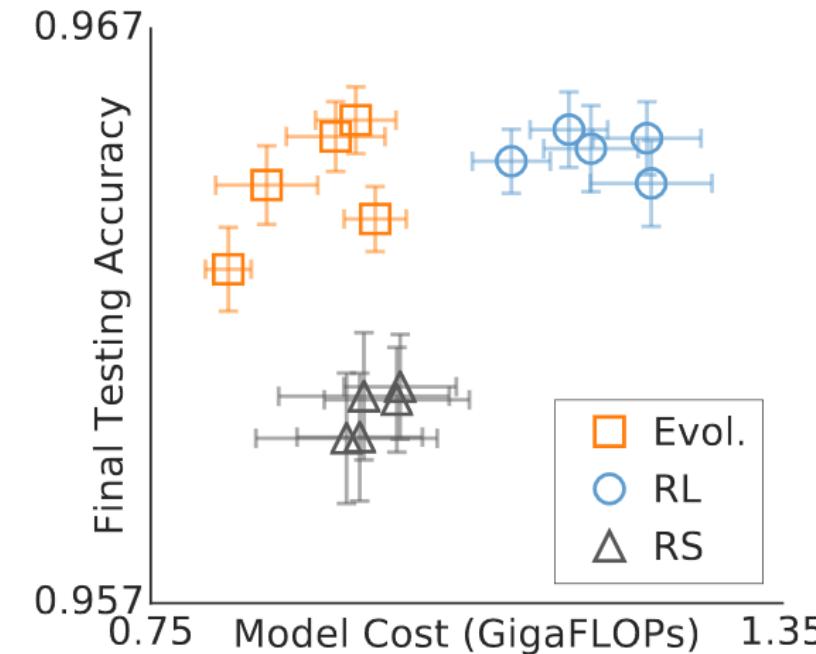


[Real et al., AAAI 2019]

during architecture search



final evaluation

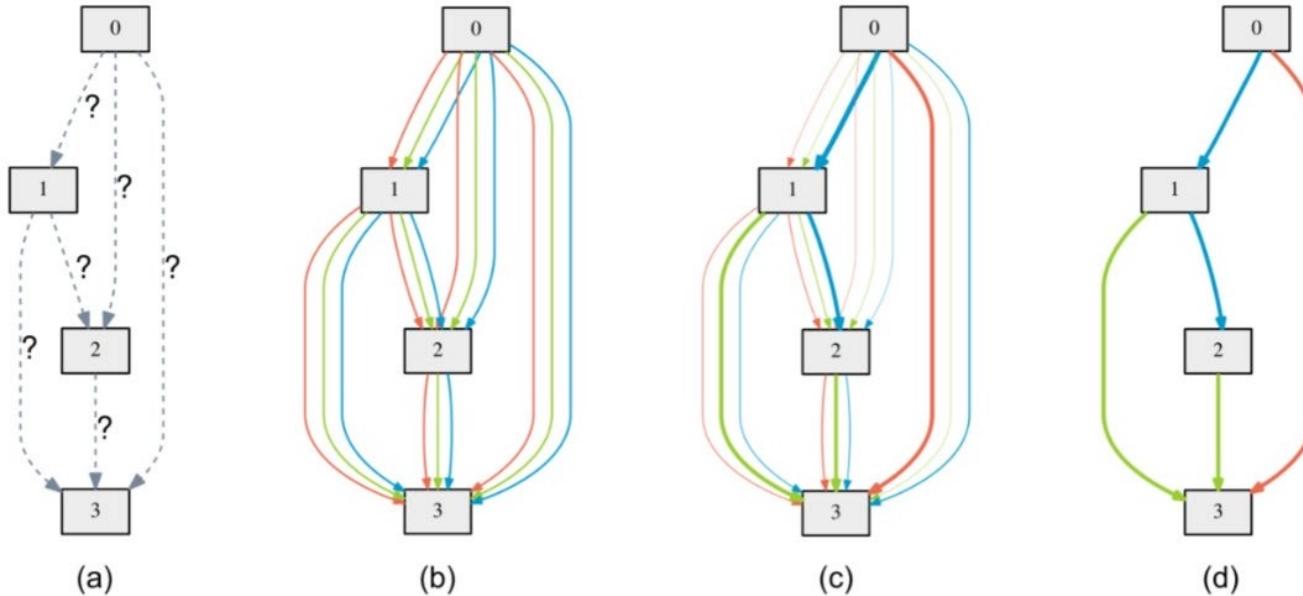


# Overview of NAS Speedup Techniques



- **Weight Sharing & One-Shot Models**
  - ENAS [Pham et al, 2018], DARTS [Liu et al, 2019] and many follow-ups
- **Weight Inheritance & Network Morphisms**
  - Local changes in architecture, followed by fine-tuning steps
  - [Cai et al, 2018; Elsken et al, 2017; Cortes et al, 2017; Cai et al, 2018, Elsken et al, 2019]
- **Meta-Learning**
  - Learning across datasets
  - To initialize architectural weights of DARTS [Lian et al, 2020; Elsken et al, 2020]
  - Prior for blackbox optimization methods [Wong et al, 2018; Runge et al, 2019; Zimmer et al, 2020]
- **Multi-Fidelity Optimization**
  - Exploit cheaper proxy models for blackbox optimizers, in particular Bayesian optimization
  - [Jamieson & Talwalkar, 2016; Li et al, 2017; Falkner et al, 2018; Zela et al, 2018; White et al, 2021]

# DARTS: Differentiable Architecture Search

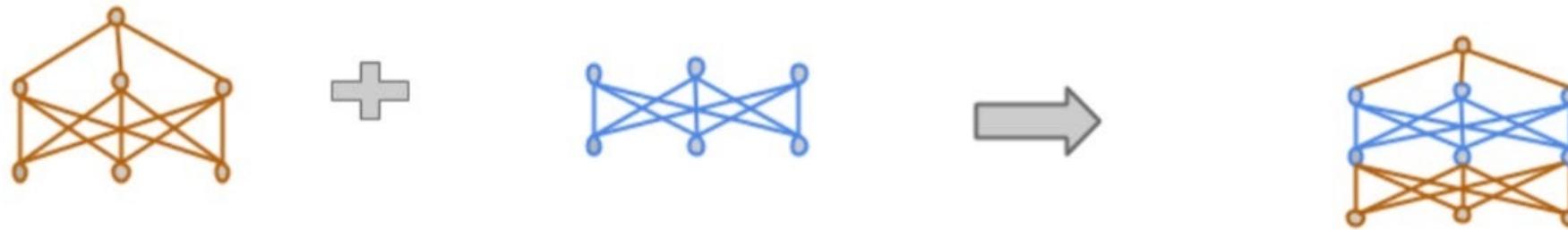


- Relax the discrete NAS problem (a->b)
  - One-shot model with continuous architecture weight  $\alpha$  for each operator
  - Mixed operator:  $\bar{o}^{(i,j)}(x) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})} o(x)$
- Solve a bi-level optimization problem (c)
$$\begin{aligned} & \min_{\alpha} \quad \mathcal{L}_{val}(w^*(\alpha), \alpha) \\ \text{s.t. } & w^*(\alpha) = \operatorname{argmin}_w \mathcal{L}_{train}(w, \alpha) \end{aligned}$$
- In the end, discretize to obtain a single architecture (d)

# Network Morphisms



- Network morphisms [Chen et al., 2016; Wei et al., 2016]
  - Change the network structure, but not the modelled function (i.e., for every input, the network yields the same output)
  - As before applying the network morphism)



- Can use this in NAS algorithms as operations to generate new networks
- Avoids costly training from scratch

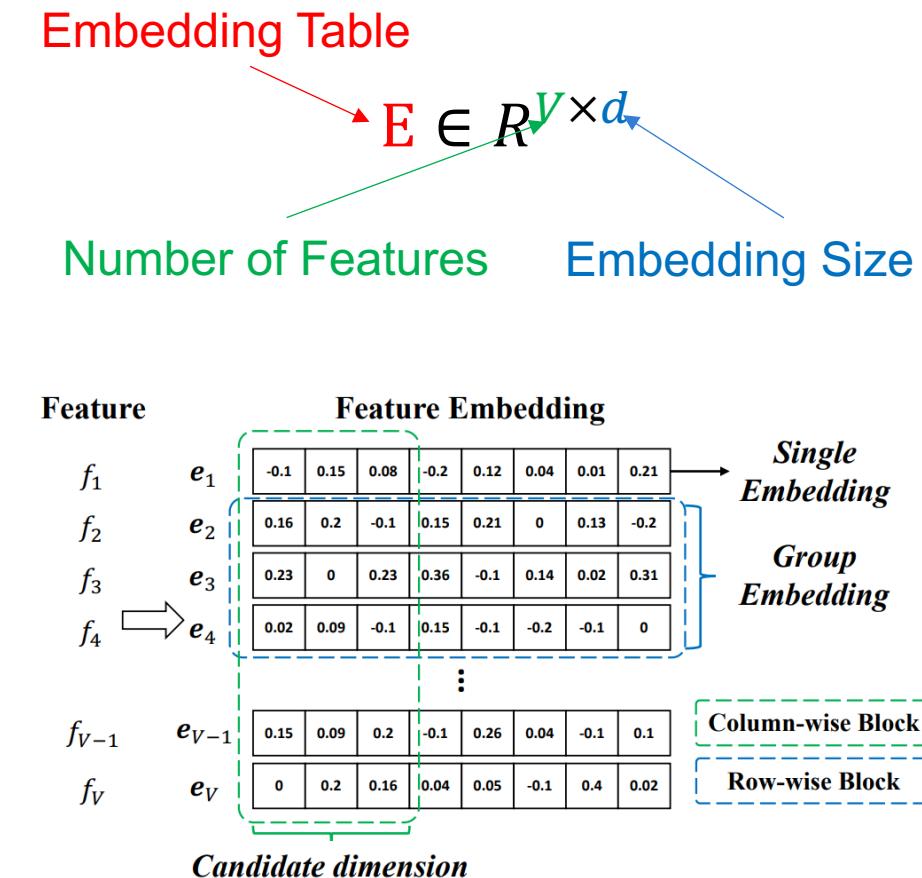


# Table of Contents

- Introduction
  - Background: Deep Recommender Systems
- Preliminary of AutoML
- **DRS Embedding Components**
  - Single Embedding Search
  - Group Embedding Search
- DRS Interaction Components
  - Feature Interaction Search
  - Interaction Function Search
  - Interaction Block Search
- DRS Comprehensive Search & System
- Conclusion & Future Direction

# DRS Embedding Components

- The Embedding Table is used to map **the high-dimensional features** into a **low-dimensional latent space**.
- Feature embedding is **the cornerstone of the DRS** as the number of parameters in DRS is concentrated in the embedding table.
- To improve the **prediction accuracy, save storage space and reduce model capacity**, AutoML-based solutions are proposed for the learning of embedding table.
  - Single Embedding Search
  - Group Embedding Search



# Single Embedding Search-AMTL



- AMTL leverages a twins-based architecture to avoid the unbalanced parameters update problem due to the different frequencies.
- AMTL designs an embedding search space with  $d^V$  size, where  $d$  is the embedding size.
- The twins-based architecture acts as a frequency-aware policy network to search the optimum dimension for each feature value → relaxed to a continuous space by tempered softmax.

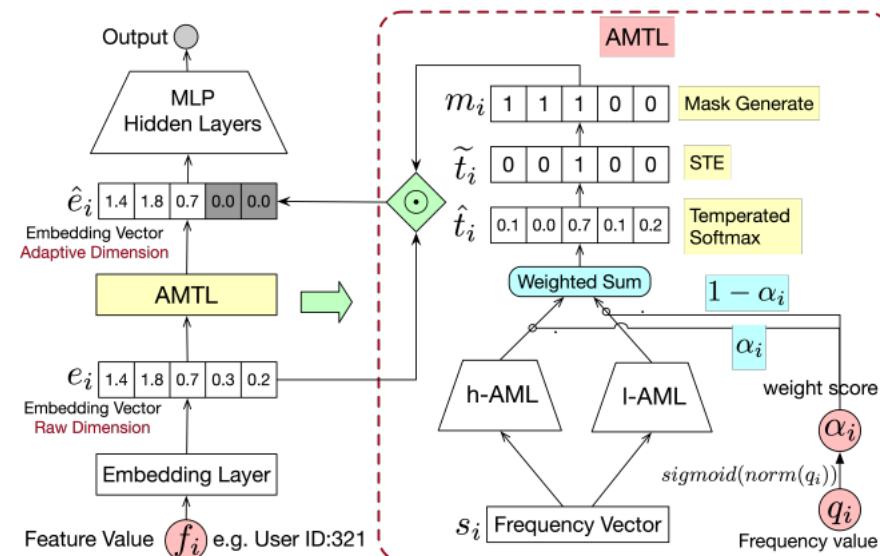


Figure 2: The framework of AMTL.

# Single Embedding Search-PEP



- Pruning-based Solution by enforcing column-wise sparsity on the Embedding Table with  $L_0$  normalization.
- Search Space:  $2^{Vd}$ , where  $V$  is number of features and  $d$  is the embedding size for each feature.

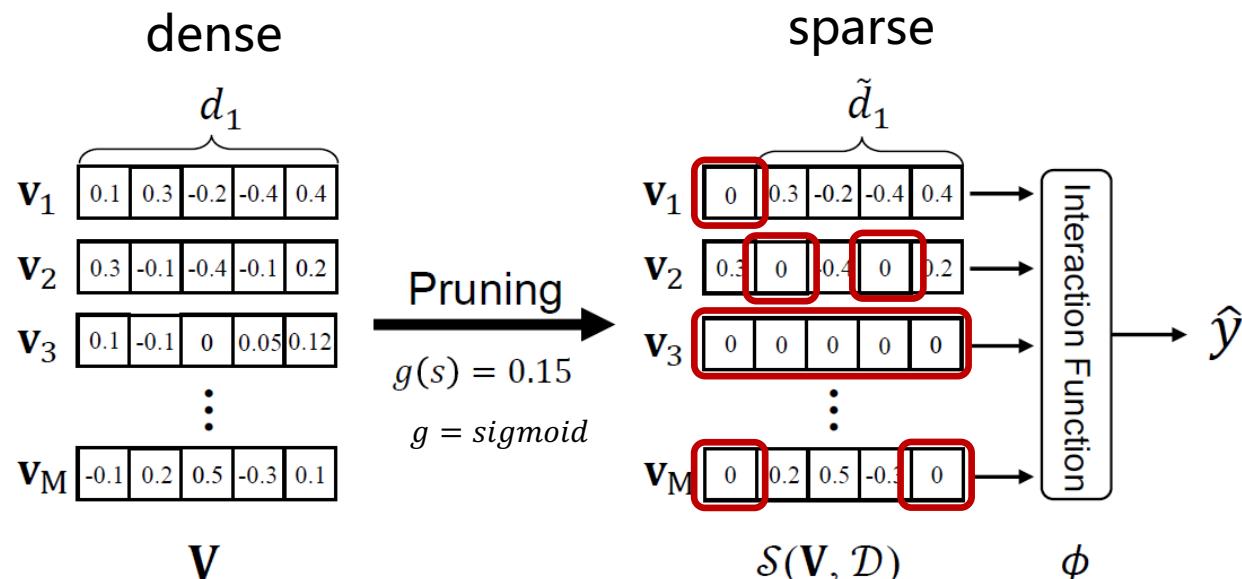


Figure 1: The basic idea of PEP.

$\min \mathcal{L}, \text{s.t. } \|\mathbf{V}\|_0 \leq k,$  NP-hard

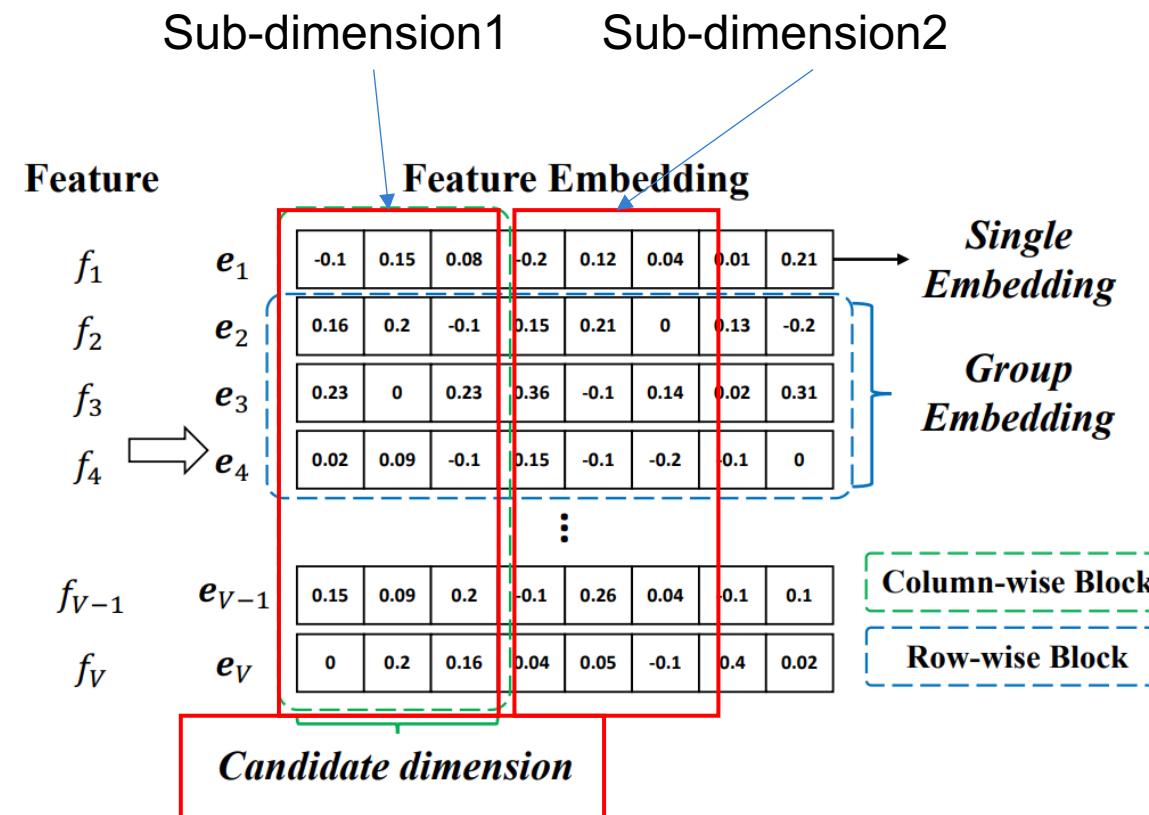
Soft threshold re-parameterization :

$$\hat{\mathbf{V}} = \mathcal{S}(\mathbf{V}, s) = \text{sign}(\mathbf{V}) \text{ReLU}(|\mathbf{V}| - g(s)),$$
$$\min \mathcal{L}(\mathcal{S}(\mathbf{V}, s), \Theta, \mathcal{D}).$$

# Single Embedding Search



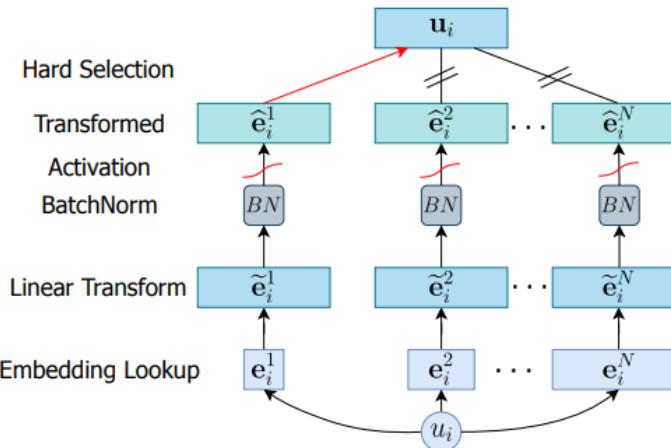
- The search space of PEP and AMTL is highly related with the embedding size d.
- To reduce the search space, AutoEmb and ESPAN divide the embedding dimension into several columnwise sub-dimensions.



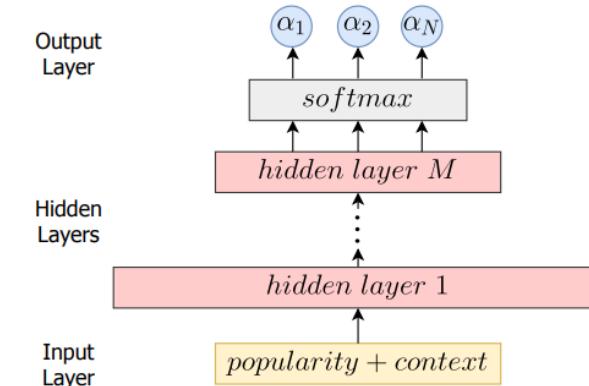
# Single Embedding Search-AutoEmb



- AutoEmb leverages two controller networks to decide the embedding sizes for users and items via end-to-end differentiable soft selection.
- Summing over the candidate subdimensions with learnable weights.
- Search Space: from  $d^V$  to  $a^V$ , where  $V$  is number of features and,  $d$  is the embedding size, and  $a$  is the number of sub-dimensions for each feature.



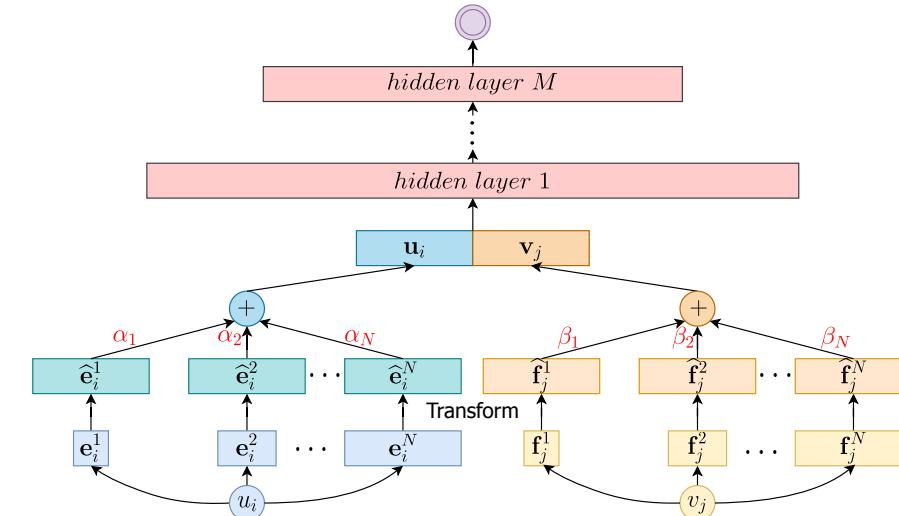
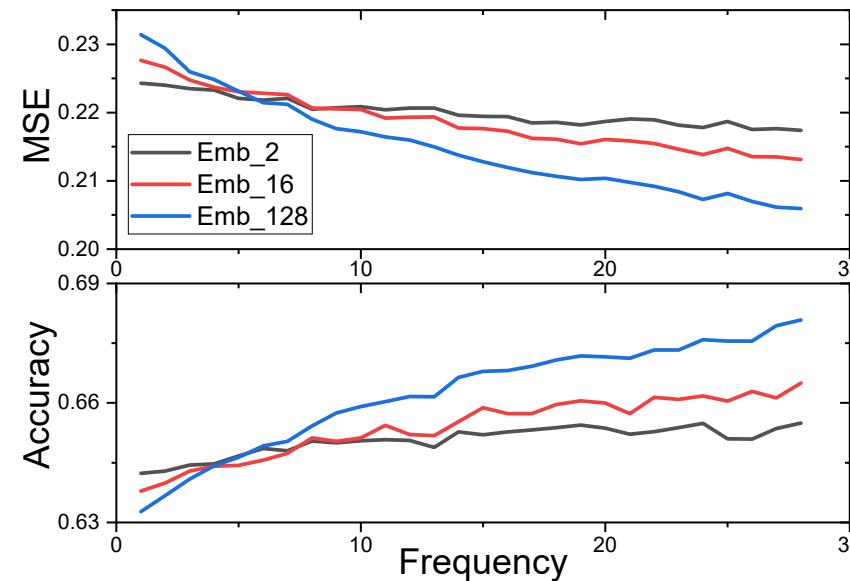
AutoEmb



Controller network

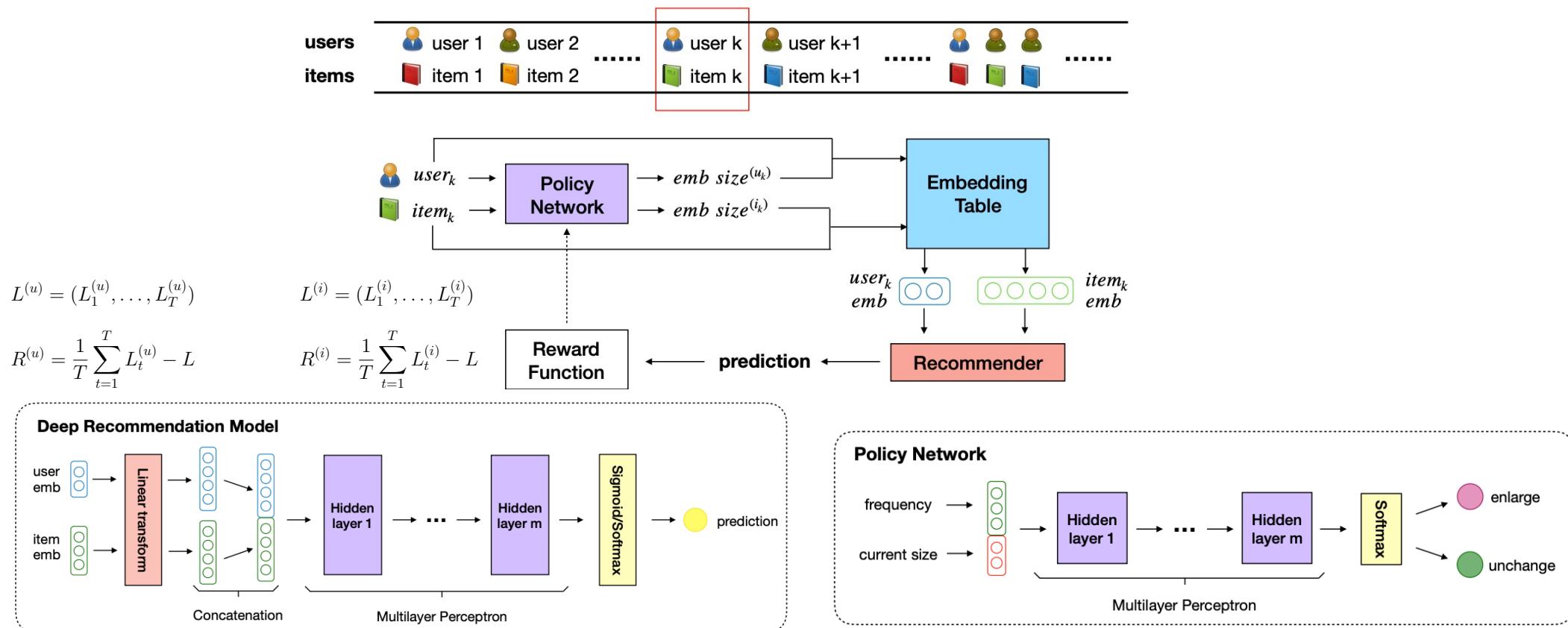
# Single Embedding Search-ESPN

- Embedding dimension often determines the capacity to encode information.
- Dynamically search the embedding sizes for different users and items
  - Optimal recommendation quality all the time
  - More efficient in memory



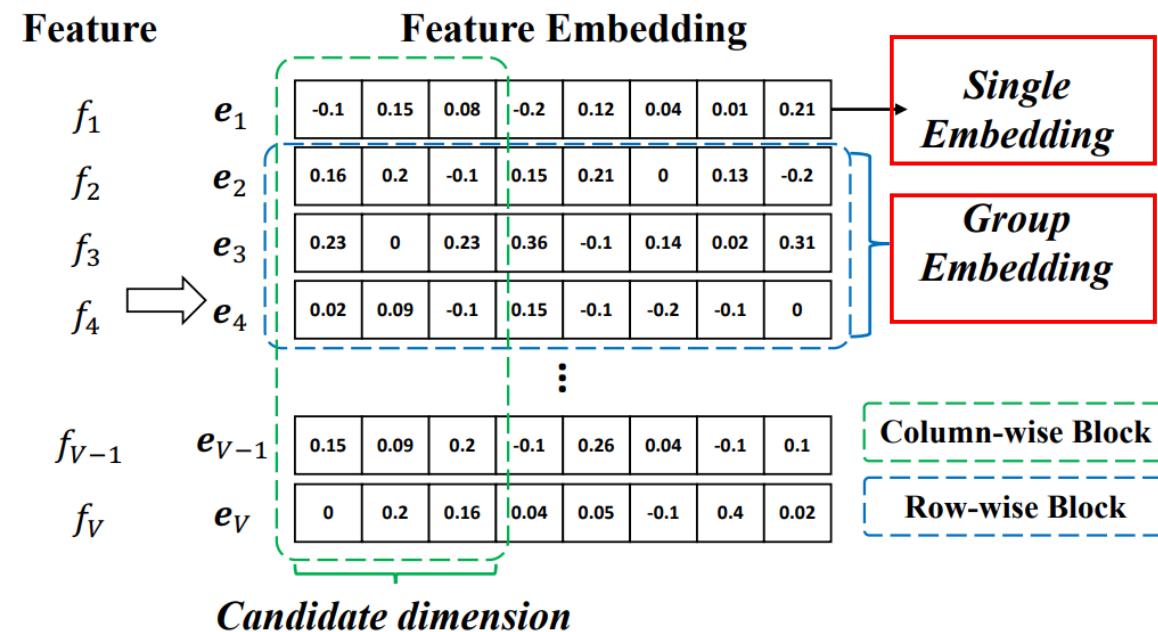
# Single Embedding Search-ESPN

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



# Group Embedding Search

- AutoEmb and ESAPN shrink the search space by dividing the embedding dimension into candidate column-wise sub-dimensions.
- Another solution is to **group** the feature values of a field based on some indicators (e.g., frequencies) and assign a **row-wise group embedding dimension** for all the values within the group.



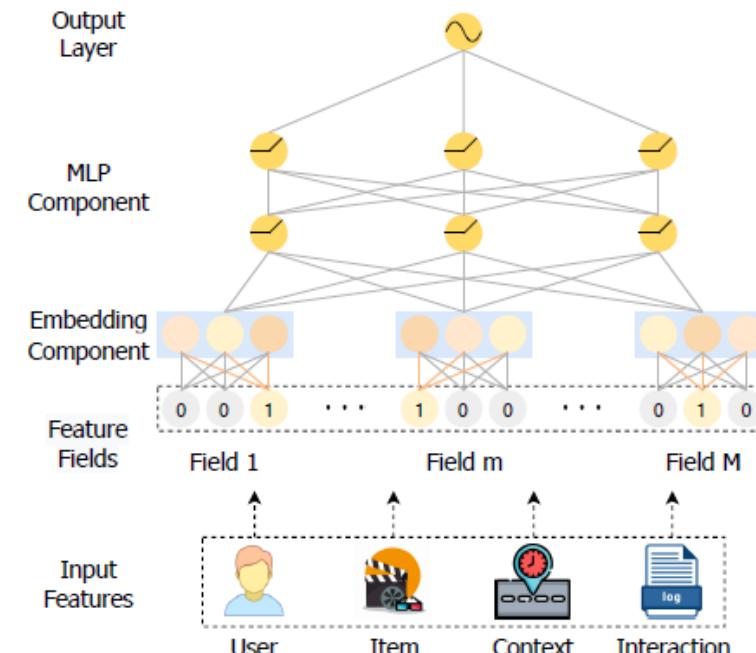
# Group Embedding Search-AutoDim



- Pre-defines several candidate sub-dimensions like AutoEmb.
- Setting the number of groups  $b = 1$  and searching a global embedding dimension for all the feature values of a **field**.
- Search Space: from  $a^V$  to  $a^m$ , where  $m$  is the number of groups.

Goal:

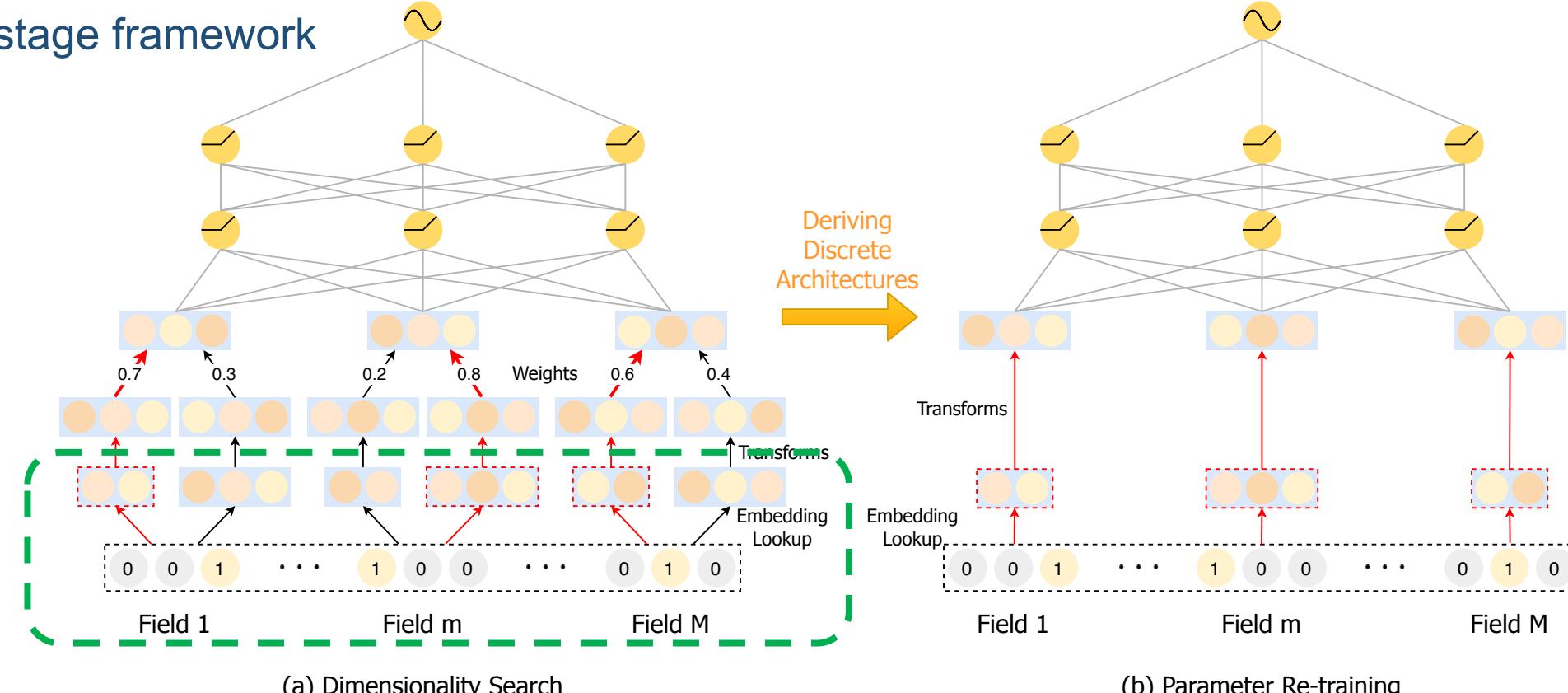
Selecting embedding dimensions to different feature fields automatically in a data-driven manner.



# Group Embedding Search-AutoDim

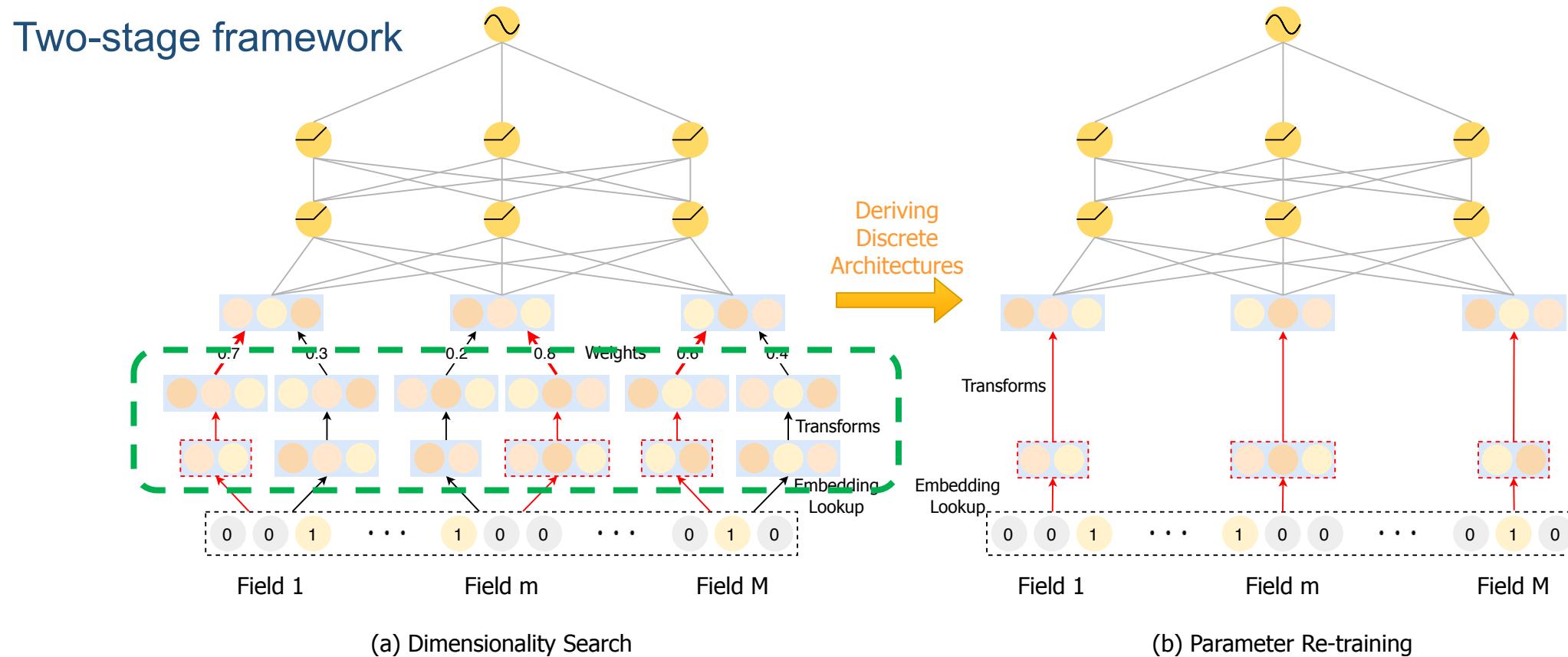


Two-stage framework



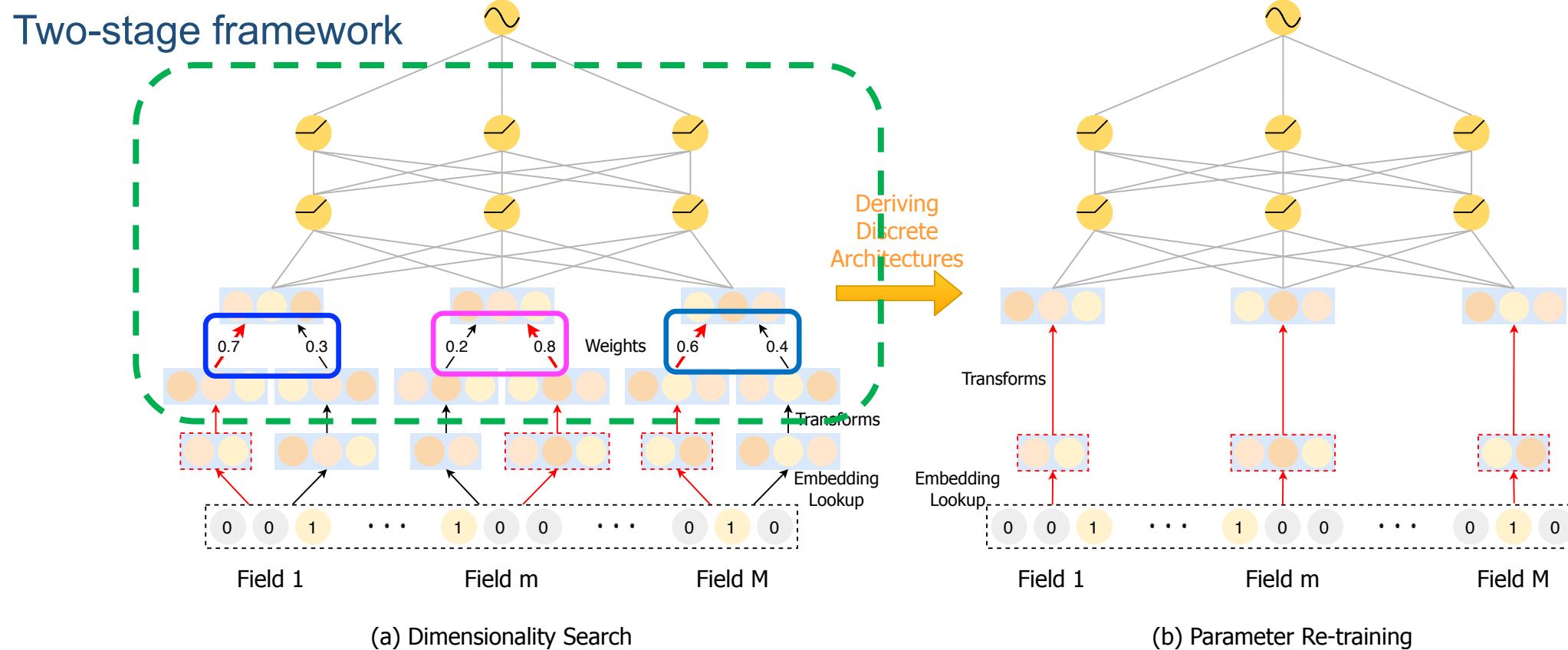
AutoDim searches dimensions for feature fields in a soft and continuous fashion via **Gumbel Softmax**, reducing to a smaller search space: 5 candidate for each feature field.

# Group Embedding Search-AutoDim



AutoDim searches dimensions for feature fields in a soft and continuous fashion via **Gumbel Softmax**, reducing to a smaller serach space: 5 candidate for each feature field.

# Group Embedding Search-AutoDim

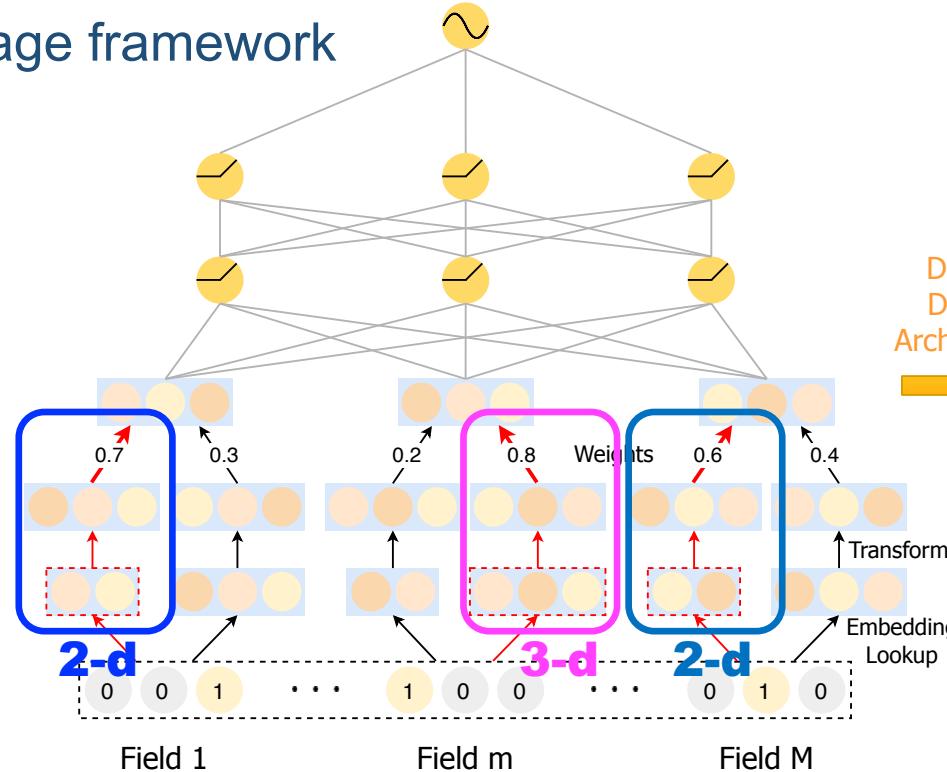


AutoDim searches dimensions for feature fields in a soft and continuous fashion via **Gumbel Softmax**, reducing to a smaller search space: 5 candidate for each feature field.

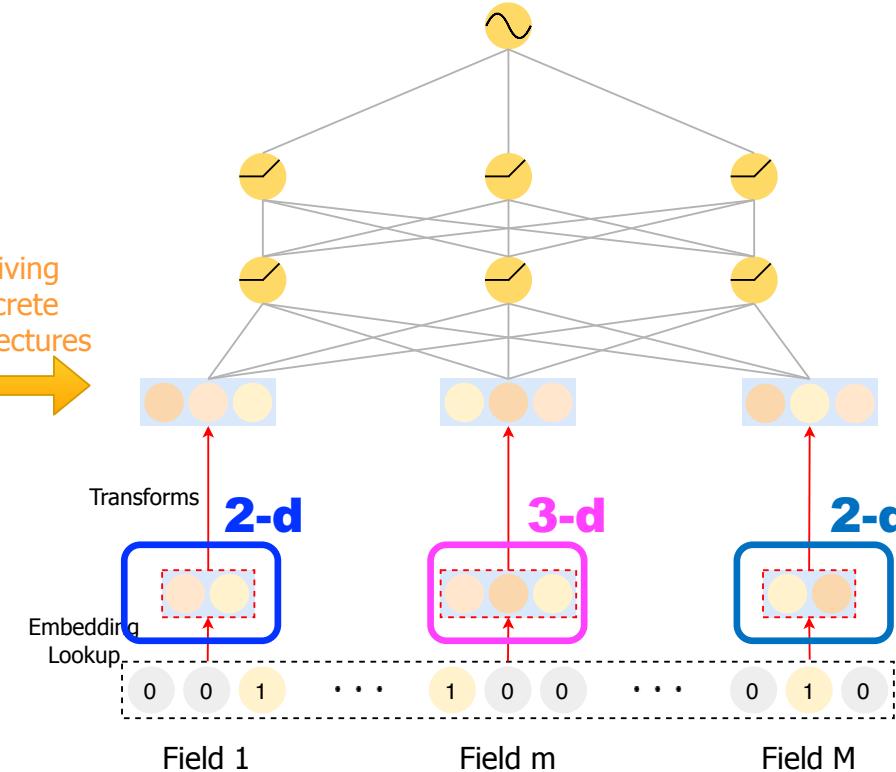
# Group Embedding Search-AutoDim



# Two-stage framework



### (a) Dimensionality Search



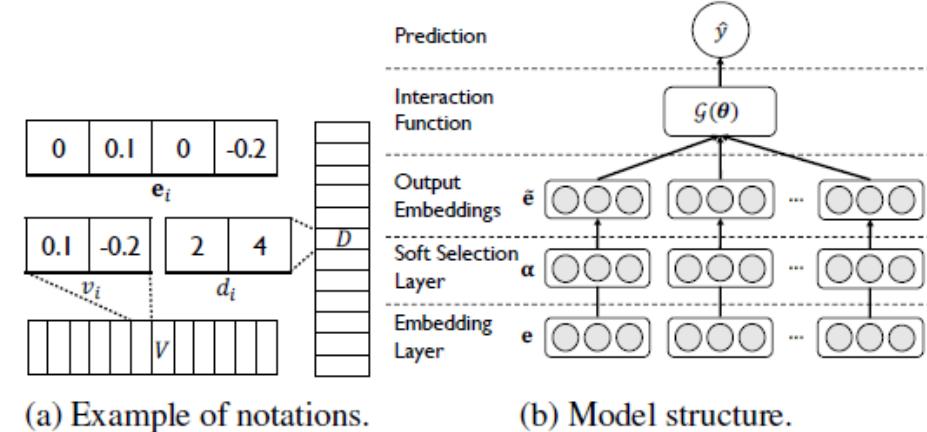
### (b) Parameter Re-training

AutoDim searches dimensions for feature fields in a soft and continuous fashion via **Gumbel Softmax**, reducing to a smaller search space: 5 candidate for each feature field.

# Group Embedding Search-DNIS



- DNIS splits the feature into multi-groups based on the feature frequencies or clustering.
- Search space: from  $2^{Vd}$  into  $2^{bd}$ , where b is the number of groups.
- DNIS searches for mixed feature embedding dimensions in a more flexible space through continuous relaxation and differentiable optimization.



(a) Example of notations.

(b) Model structure.

$$\tilde{e}_i = e_i \odot \alpha_{l_*}$$

$$\tilde{E}_{i,j} = \begin{cases} 0, & \text{if } |\tilde{E}_{i,j}| < \epsilon \\ \tilde{E}_{i,j}, & \text{otherwise} \end{cases}$$

# Group Embedding Search-NIS



- Input component assigns embedding vectors to each item of these discrete features, which dominate both the size and the inductive bias of the model.
- The vocabulary and embedding sizes for discrete features are often selected heuristically.

## Head Feature

- More data, more information
- Needing larger embedding size

## Tail Feature

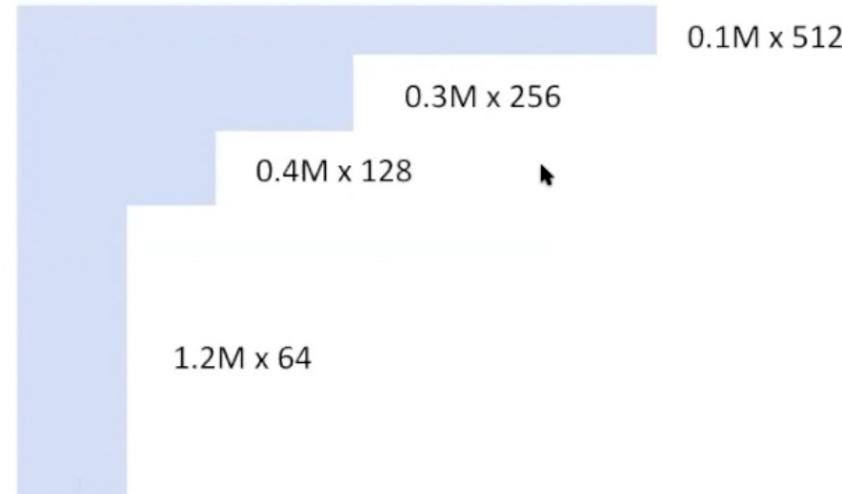
- Less data, less information
- Small embedding size is enough

Single-size Embedding (SE)

$$\sum_{F \in \mathcal{F}} v_F \times d_F \leq \mathcal{C}$$

Multi-size Embedding (ME)

$$\sum_{F \in \mathcal{F}} \sum_{i=1}^{M_F} v_{F_i} \times d_{F_i} \leq \mathcal{C}$$

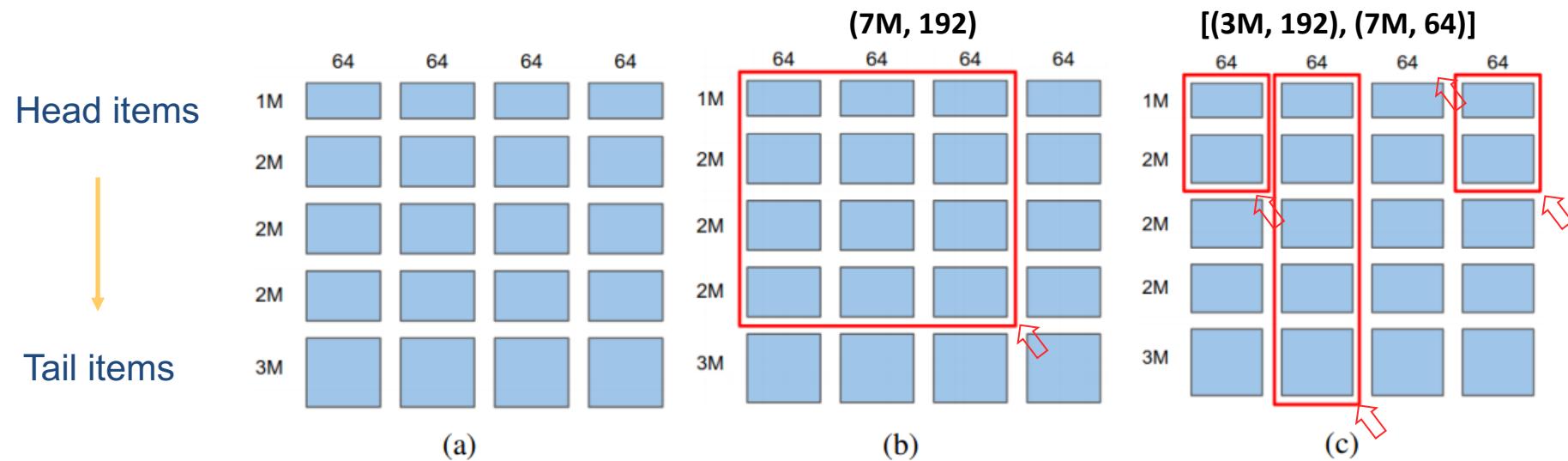
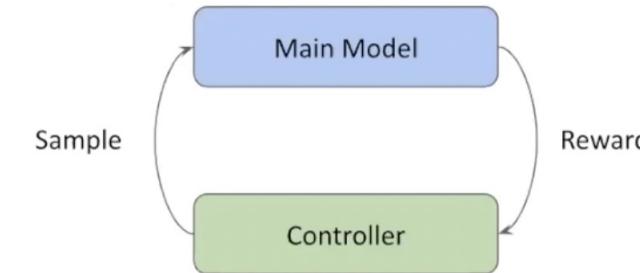


# Group Embedding Search-NIS



## RL-based AutoML approach (ENAS)

- Main model is the deep recommendation model
- Controller learns to sample embedding dimensions that generate higher reward.
- Reward:  $R = R_Q - \lambda * C_M$



**Single-size Embedding (SE)   Multi-size Embedding (ME)**

**Embedding Blocks:** discretizing an embedding matrix of size  $v \times d$  into  $S \times T$  sub-matrices

# Group Embedding Search-AutoDis



Existing methods for numerical feature representation have some limitations:

## 1. Category 1: No Embedding

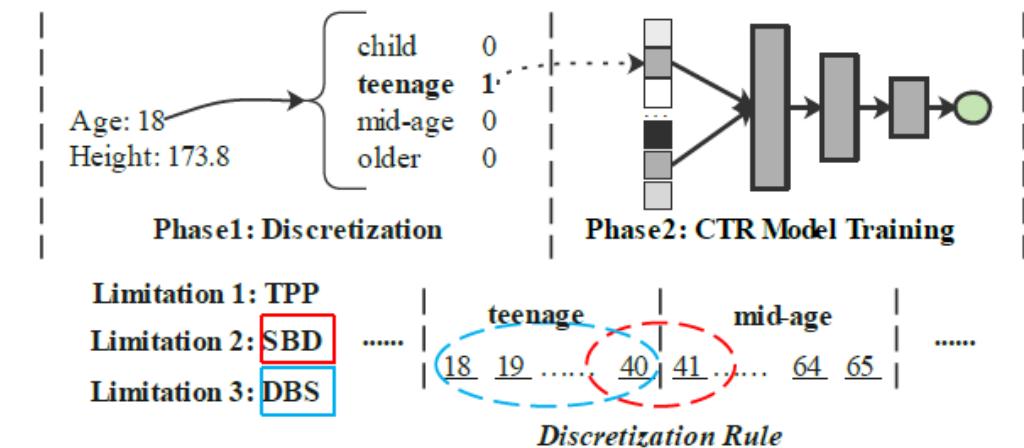
- **Low capacity**: difficult to capture informative knowledge of numerical fields.
- **Poor compatibility**: difficult to adapt to some models (e.g., FM).

## 2. Category 2: Field Embedding

- **Low capacity**: single shared field-specific embedding.

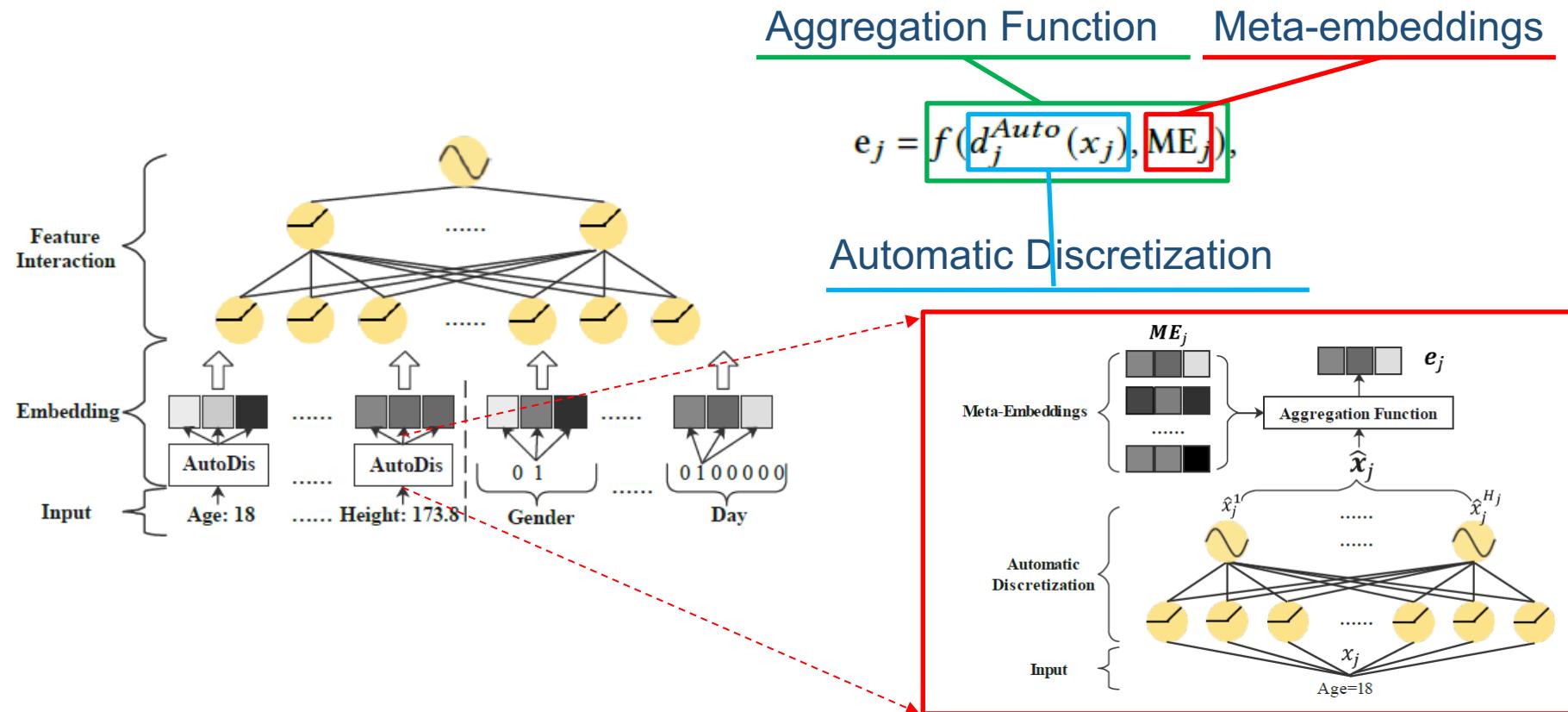
## 3. Category 3: Discretization

- **TPP** (Two-Phase Problem)
- **SBD** (Similar value But Dis-similar embedding)
- **DBS** (Dis-similar value But Same embedding)



# Group Embedding Search-AutoDis

AutoDis is a numerical features embedding learning framework with **high model capacity, end-to-end training and unique representation** properties preserved.



# Summarize DRS Embedding



	Approach	Search Space	Search Type
AMTL	DARTS	$d^V$	Single Embedding Search
PEP	L0	$2^{Vd}$	Single Embedding Search
AutoEmb	DARTS	$a^V$	Single Embedding Search
ESSPAN	RL	$a^V$	Single Embedding Search
AutoDim	DARTS	$a^b$	Group Embedding Search
NIS	DARTS	$b^a$ or $ab$	Group Embedding Search
DNIS	darts	$2^{bd}$	Group Embedding Search
AutoDis	Softmax+MLP	--	Group Embedding Search

\*We do not show the search space as it is the embedding approach for numerical feature.

\*\* d is the embedding size, b is the number of feature groups, V is the number of features, a is the candidate embedding sizes, m is ? a < d, b<<V.

- DARTS is more popular as it has higher efficiency.
- The Search Space of Group Embedding Search is less than Single Embedding Search.
- Limited approaches for embedding learning of numerical features,



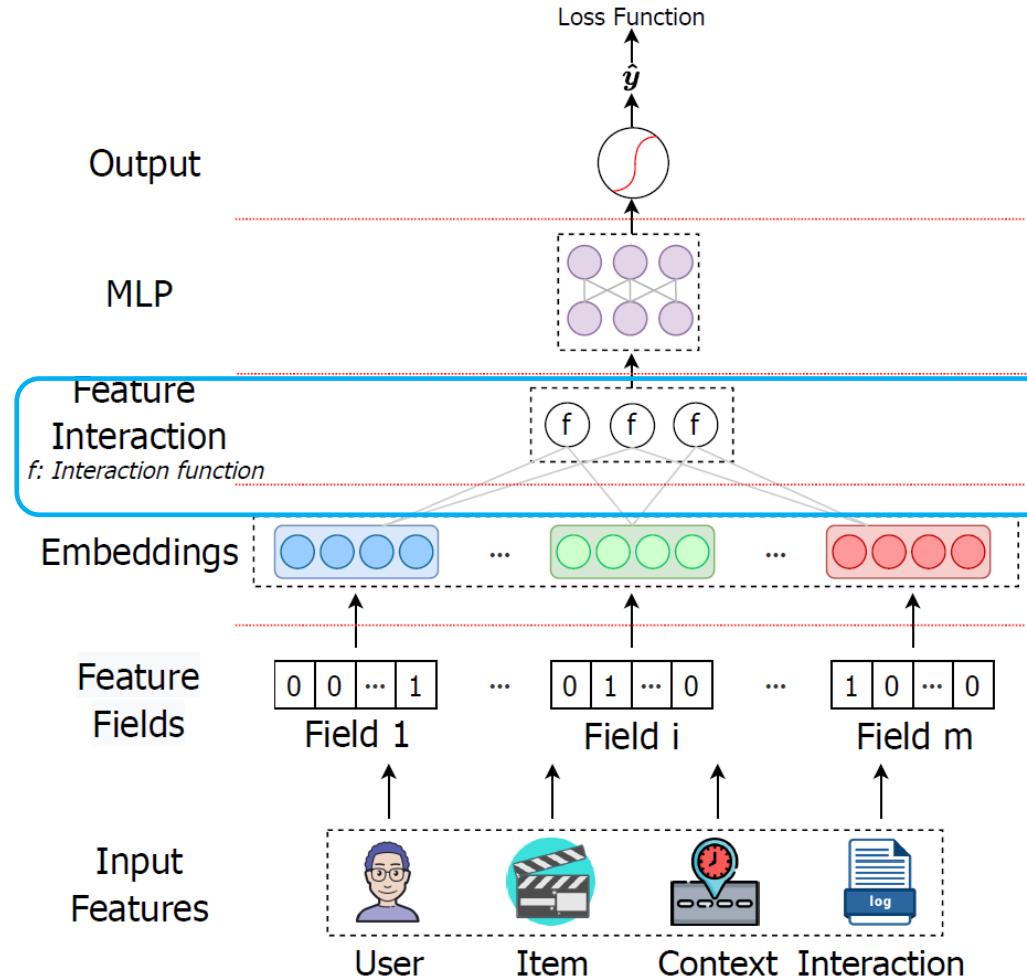
# Table of Contents

- Introduction
  - Background: Deep Recommender Systems
- Preliminary of AutoML
- DRS Embedding Components
  - Single Embedding Search
  - Group Embedding Search
- **DRS Interaction Components**
  - Feature Interaction Search
  - Interaction Function Search
  - Interaction Block Search
- DRS Comprehensive Search & System
- Conclusion & Future Direction

# Background



Effectively modelling **feature interactions** is important.



- Both low-order and high-order feature interactions play important roles to model user preference.
  - 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 → 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)

The challenges of modelling **feature interactions**:

1) Enumerate all feature interactions

- Large memory and computation cost
- Difficult to be extended into high-order interactions
- Useless interaction

2) Require human efforts to identify important feature interactions

- High labor cost
- Risks missing some counterintuitive (but important) interactions

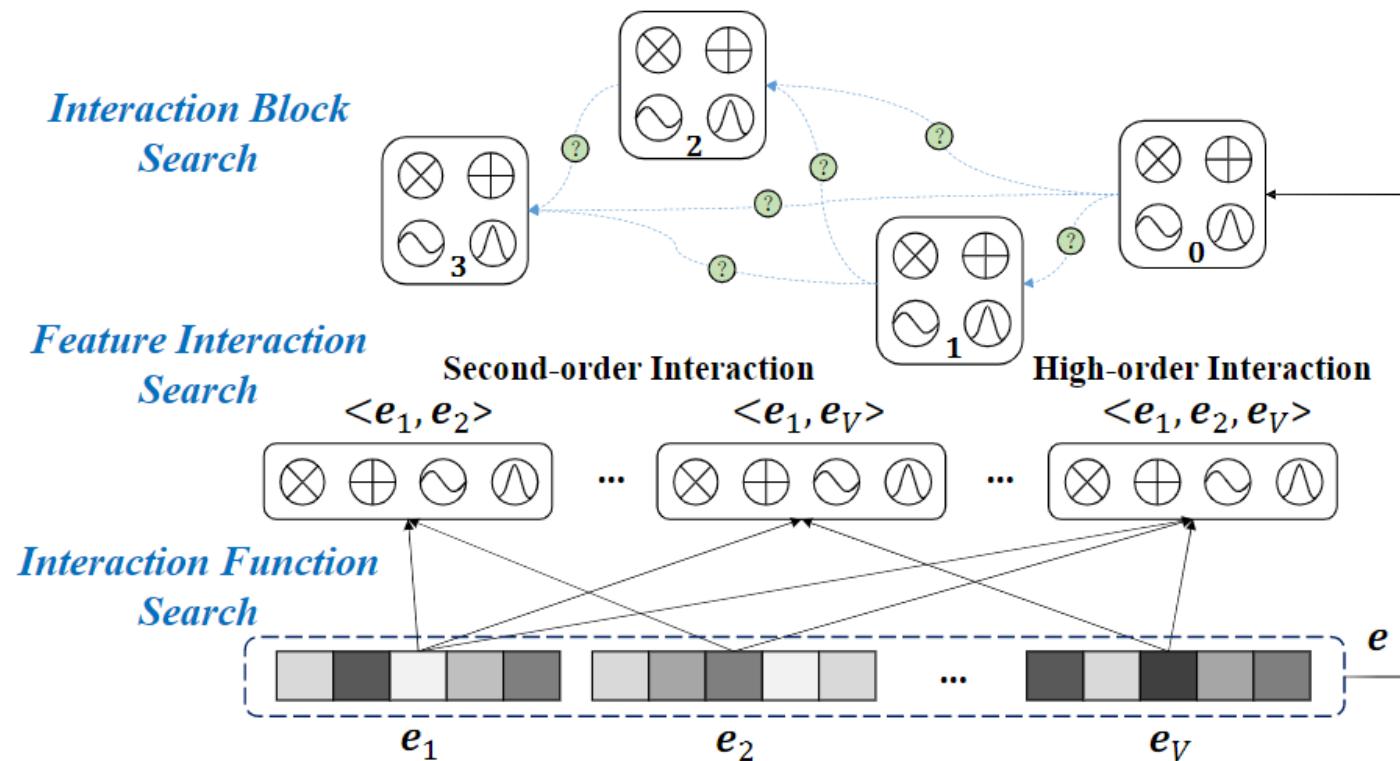
3) Require human efforts to select appropriate interaction functions

- Human expert knowledge
- Global interaction function for all the feature interactions

# Background



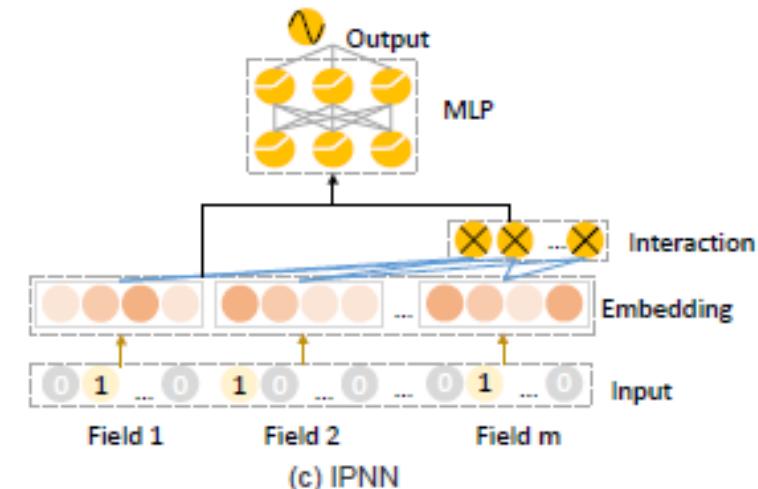
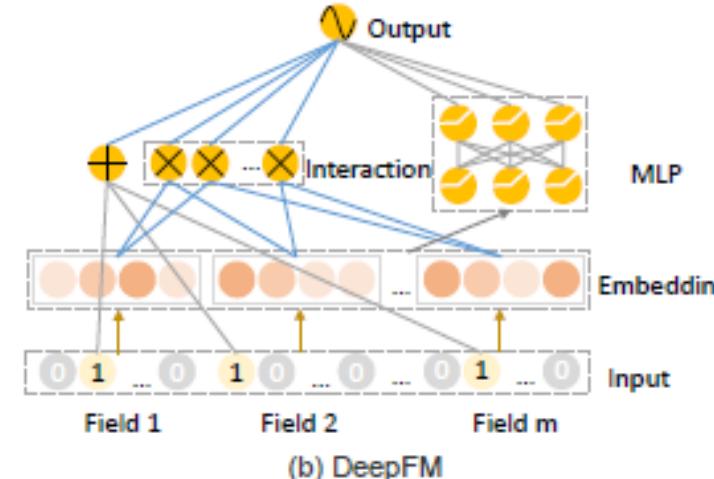
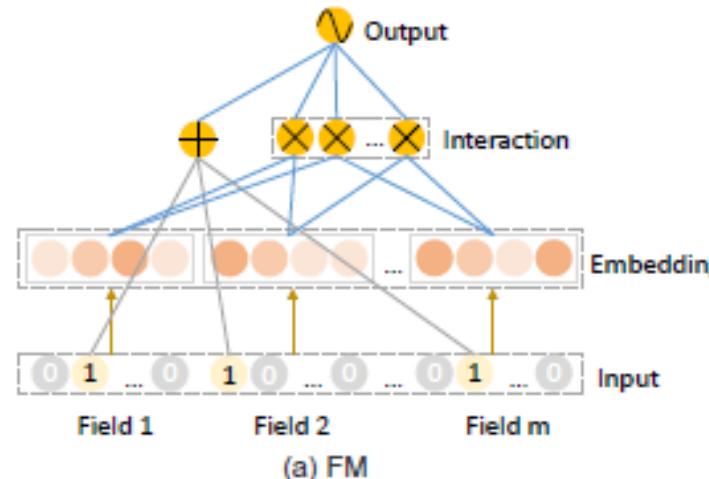
Automatically select important feature interactions with appropriate interaction functions



# Feature Interaction Search-AutoFIS



- Not all the feature interactions are useful.
- We apply AutoML techniques to identify such noisy feature interactions and filter them.
- Attention is not enough, as the learning of useful feature interactions is still diluted by many useless ones.



# Feature Interaction Search-AutoFIS



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

# Feature Interaction Search-AutoFIS

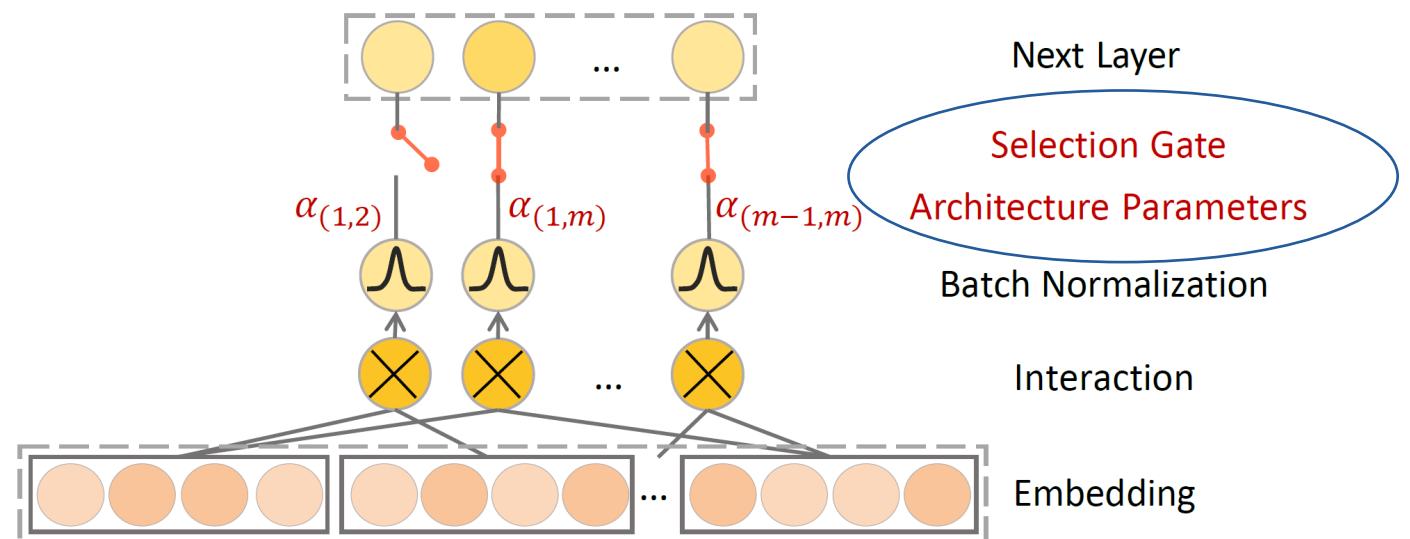


## Search Stage:

- Gate for each feature interaction
  - Huge search space  $2^{C_m^2}$
- To make such identification differentiable, we relax the discrete search space to be continuous, by defining architecture parameters  $\alpha$ .
- We use Batch Normalization and GRDA Optimizer to train stable and sparse architecture parameters.

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

Indicator  $\alpha = 0$  or  $1$



## Retrain Stage:

- Abandon unimportant feature interactions
- Retrain model

## The limitation of AutoFIS:

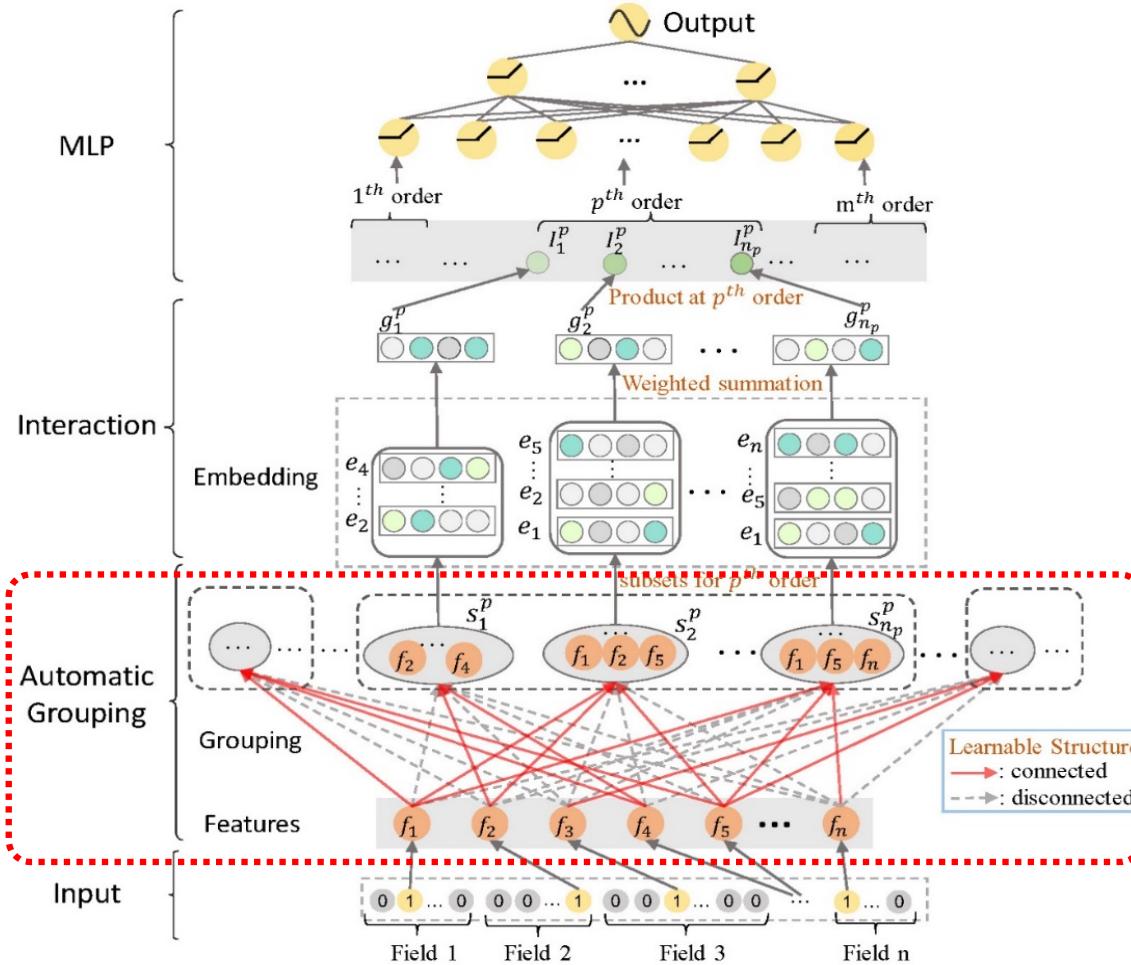
- When searching **high-order feature interactions**, the search space of AutoFIS is huge, resulting in low search efficiency.

## Solution of AutoGroup:

- To solve the **efficiency-accuracy dilemma**, AutoGroup proposes automatic feature grouping, reducing the  $p$ th-order search space from  $2^{C_m^p}$  to  $2^{gm}$ , where  $g$  is the number of pre-defined groups.

# Feature Interaction Search-AutoGroup

## 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:

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

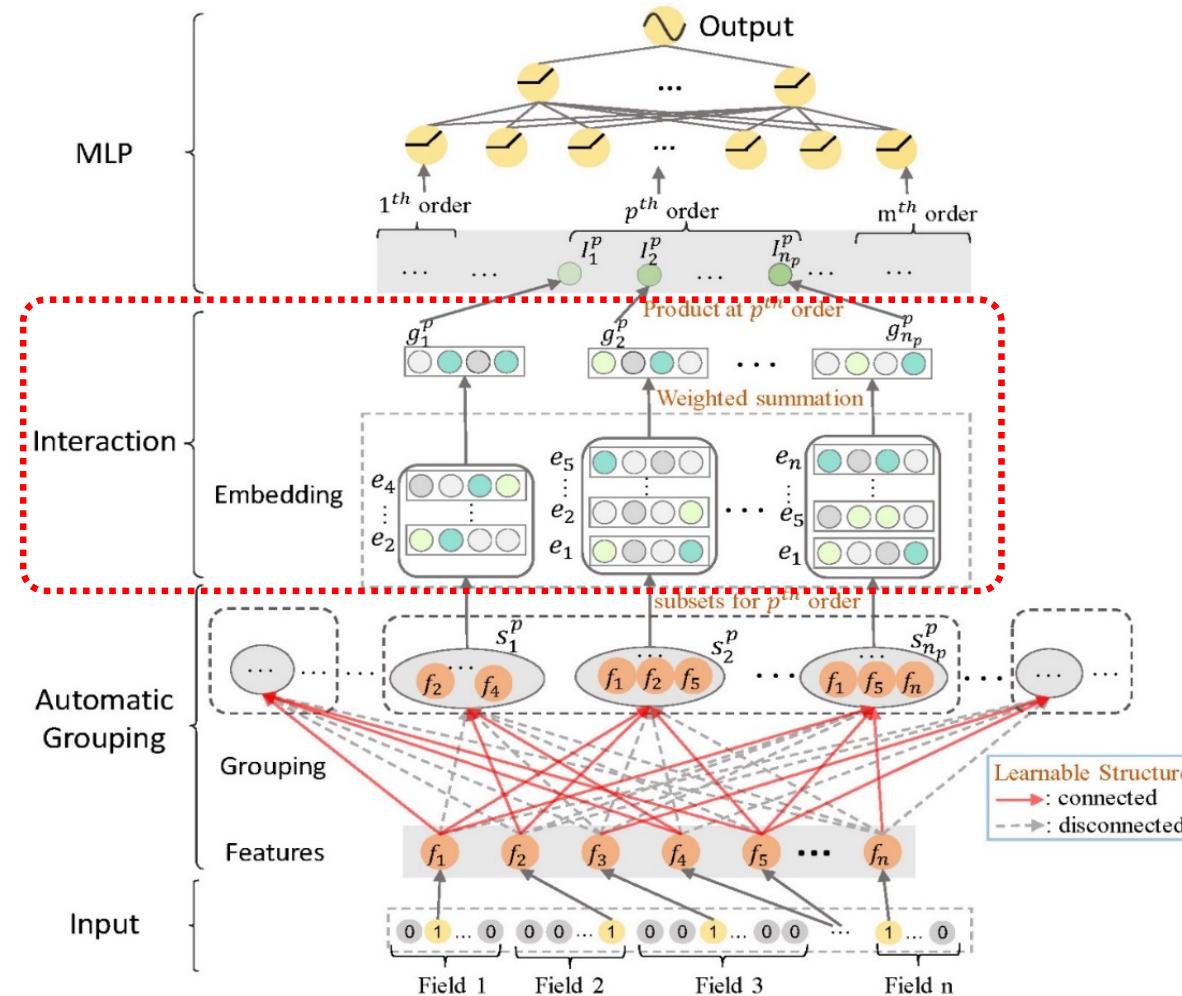
$$\alpha_0 = \frac{1}{1 + \exp(-\alpha_{i,j}^p)} \quad \alpha_1 = \frac{\exp(-\alpha_{i,j}^p)}{1 + \exp(-\alpha_{i,j}^p)}$$

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

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

# Feature Interaction Search-AutoGroup

## Interaction Stage:



Feature set representation:

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

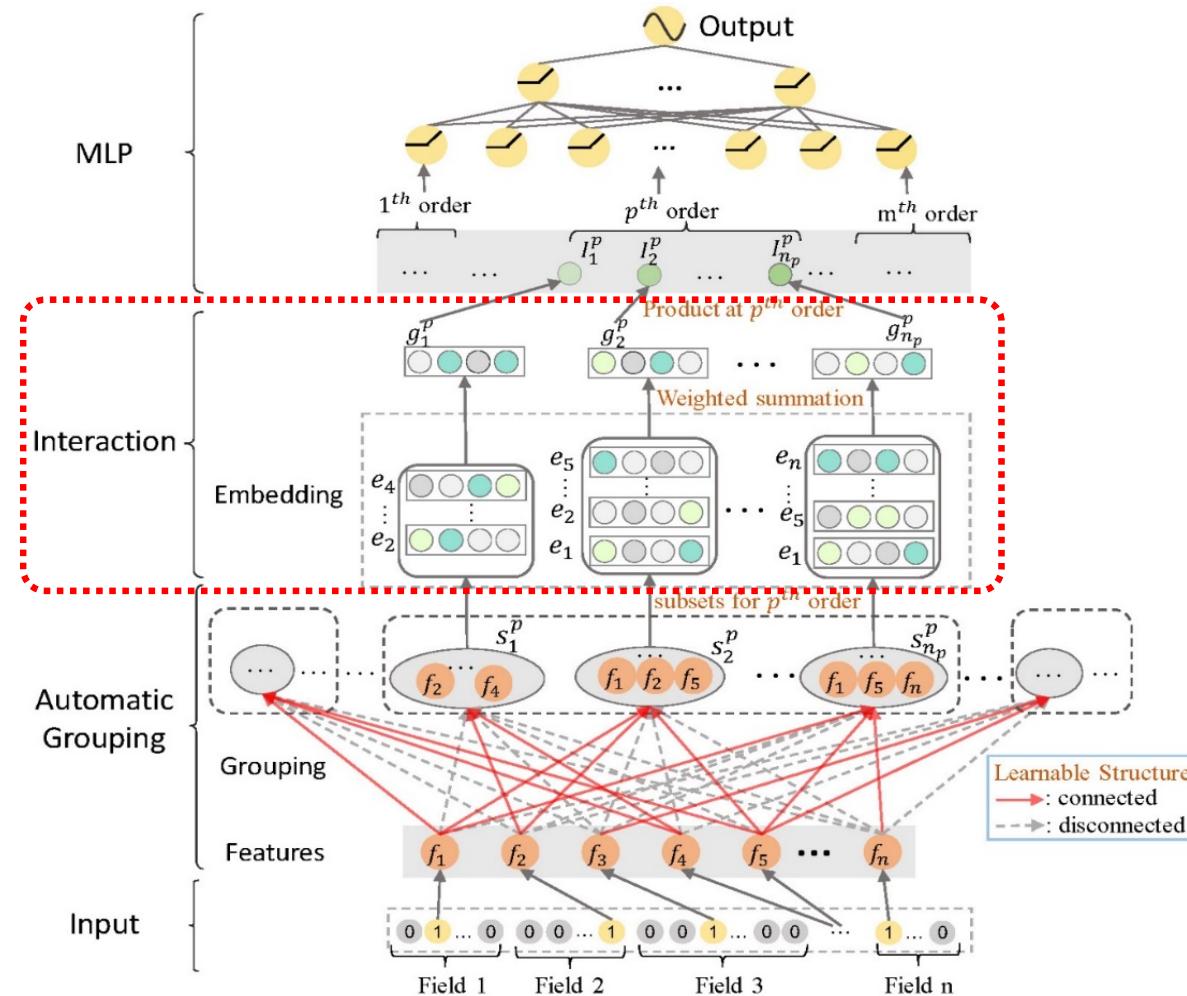
$s_j^p$ : the  $j^{th}$  feature set for order- $p$  feature interactions.

$e_i$ : embedding for feature  $f_i$

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

# Feature Interaction Search-AutoGroup

## Interaction Stage:



### Interaction at a given order:

- Inspired by the reformulation of FM:

$$\begin{aligned}\hat{y}(x) &= w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle e_i, e_j \rangle x_i x_j \\ &= w_0 + \sum_{i=1}^n w_i x_i + \frac{1}{2} \left( \left( \sum_{i=1}^n x_i e_i \right)^2 - \sum_{i=1}^n (x_i e_i)^2 \right)\end{aligned}$$

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

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

## The limitation of AutoGroup:

- Although AutoGroup solve the ***efficiency-accuracy dilemma*** via feature grouping, it ignores the ***order-priority property***.
- ***Order-priority property***: The higher-order feature interactions quality can be relevant to their de-generated low-order ones, and lower-order feature interactions are likely to be more vital compared with higher-order ones.

## Solution of FIVES:

- To solve the ***efficiency-accuracy dilemma***, FIVES regards the original features as a feature graph conceptually and models the high-order feature interactions by a GNN with layer-wise adjacency matrix, reducing the  $p$ th-order search space from  $2^{C_m^p}$  to  $2^{m^2}$ .
- To keep the ***order-priority property***, FIVES parameterizes the adjacency matrix and makes them depend on the previous layer.

# Feature Interaction Search-FIVES



- With an adjacency tensor  $A$ , the dedicated graph convolutional operator produces the node representations layer-by-layer.. For the step  $k$ :

$$n_i^{(k)} = p_i^{(k)} \odot n_i^{(k-1)}$$

(1)

where  $p_i^{(k)} = \text{MEAN}_{j|A_{i,j}^{(k)}=1} \{W_j n_j^{(0)}\}$

- The node representation at  $k$ -th layer corresponds to the generated  $(k + 1)$ -order interactive features:

$$n_i^{(k)} = \text{MEAN}_{j|A_{i,j}^{(k)}=1} \{W_j n_j^{(0)}\} \odot n_i^{(0)}$$

(2)

$$\approx \text{MEAN}_{(c_1, \dots, c_k) | A_{i,c_j}^{(j)}=1, j=1, \dots, k} \{f_{c_1} \otimes \dots \otimes f_{c_k} \otimes f_i\}$$

# Feature Interaction Search-FIVES



- The task of generating useful interactive features is equivalent to learning an optimal adjacency tensor  $A$ , namely **edge search**.
- The edge search task could be formulated as a bi-level optimization problem:

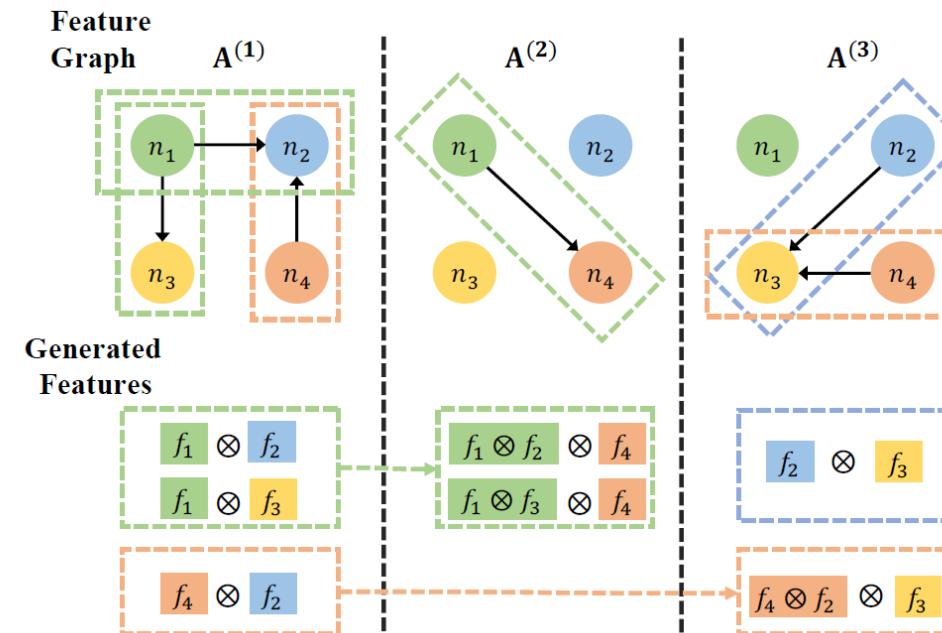
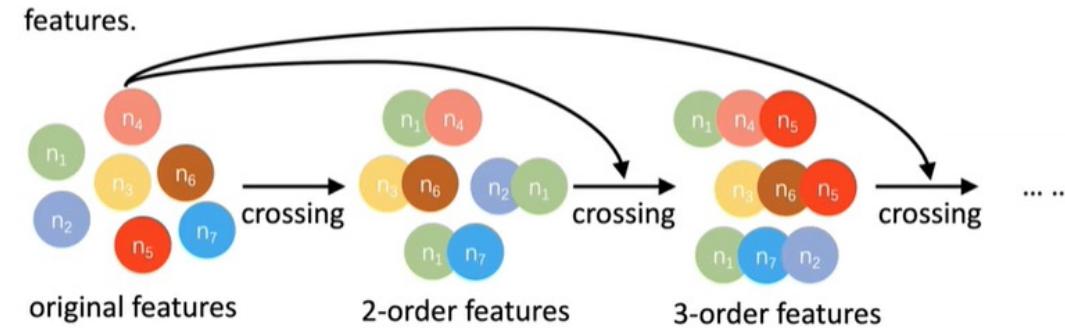
$$\begin{aligned} & \min_{\mathbf{A}} \mathcal{L}(\mathcal{D}_{\text{val}} | \mathbf{A}, \Theta(\mathbf{A})) \\ \text{s. t. } & \Theta(\mathbf{A}) = \arg \min_{\Theta} \mathcal{L}(\mathcal{D}_{\text{train}} | \mathbf{A}, \Theta) \end{aligned} \quad (3)$$

- To make the optimization more efficient, we allow to use a soft  $A^{(k)}$  for propagation at the  $k$ -th layer, while the calculation of  $A^{(k)}$  still depends on a binarized  $A^{(k-1)}$ :

$$A^{(k)} \triangleq (D^{(k-1)})^{-1} \varphi(A^{(k-1)}) \sigma(H^{(k)}) \quad (4)$$

Degree matrix of  $A^{(k-1)}$       Binarize the soft  $A^{(k-1)}$       Interactions at  $k$ -th layer

# Feature Interaction Search-FIVES



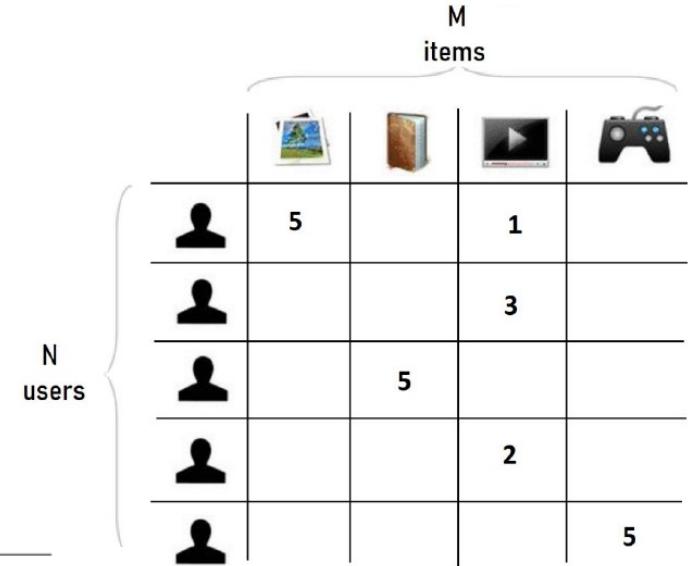
# Interaction Function Search-SIF



- Generate embedding vectors for users and items
- Generate predictions by an inner product between embedding vectors
- Evaluate predictions by a loss function on the training data set

Interaction function:  
how embedding vectors interact with each other

	IFC	operation	space	predict time	recent examples
human-designed	$\langle \mathbf{u}_i, \mathbf{v}_j \rangle$	inner product	$O((m+n)k)$	$O(k)$	MF [28], FM [37]
	$\mathbf{u}_i - \mathbf{v}_j$	plus (minus)	$O((m+n)k)$	$O(k)$	CML [19]
	$\max(\mathbf{u}_i, \mathbf{v}_j)$	max, min	$O((m+n)k)$	$O(k)$	ConvMF [25]
	$\sigma([\mathbf{u}_i; \mathbf{v}_j])$	concat	$O((m+n)k)$	$O(k)$	Deep&Wide [9]
	$\sigma(\mathbf{u}_i \odot \mathbf{v}_j + \mathbf{H}[\mathbf{u}_i; \mathbf{v}_j])$	multi, concat	$O((m+n)k)$	$O(k^2)$	NCF [17]
	$\mathbf{u}_i * \mathbf{v}_j$	conv	$O((m+n)k)$	$O(k \log(k))$	ConvMF [25]
	$\mathbf{u}_i \otimes \mathbf{v}_j$	outer product	$O((m+n)k)$	$O(k^2)$	ConvNCF [16]
AutoML	SIF (proposed)	searched	$O((m+n)k)$	$O(k)$	—



Is there an absolute best IFC? : NO, depends on tasks and datasets [1]

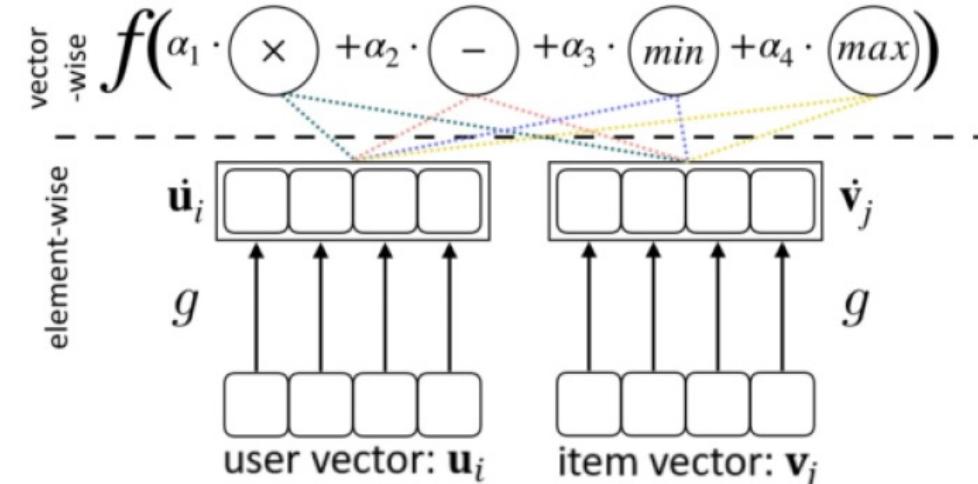
# Interaction Function Search-SIF



- SIF selects different interaction functions across different datasets.

IFC	operation
$\langle \mathbf{u}_i, \mathbf{v}_j \rangle$	inner product
$\mathbf{u}_i - \mathbf{v}_j$	plus (minus)
$\max(\mathbf{u}_i, \mathbf{v}_j)$	max, min
$\sigma([\mathbf{u}_i; \mathbf{v}_j])$	concat
$\sigma(\mathbf{u}_i \odot \mathbf{v}_j + \mathbf{H}[\mathbf{u}_i; \mathbf{v}_j])$	multi, concat
$\mathbf{u}_i * \mathbf{v}_j$	conv
$\mathbf{u}_i \otimes \mathbf{v}_j$	outer product

Cut the search space into two blocks

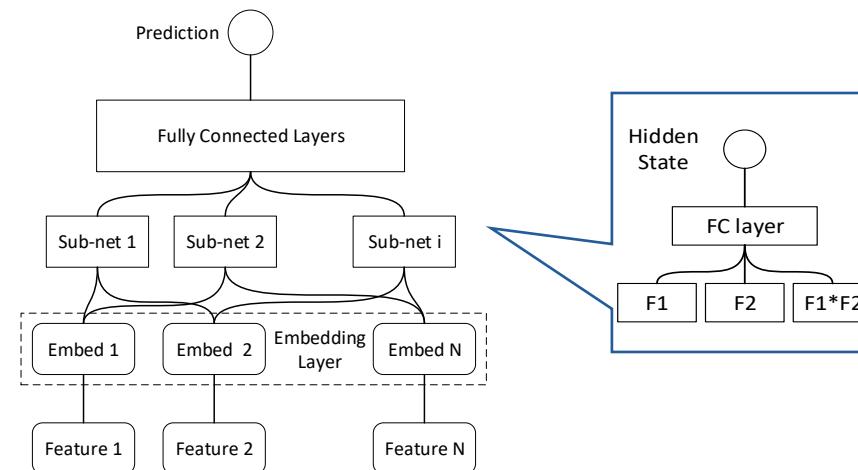
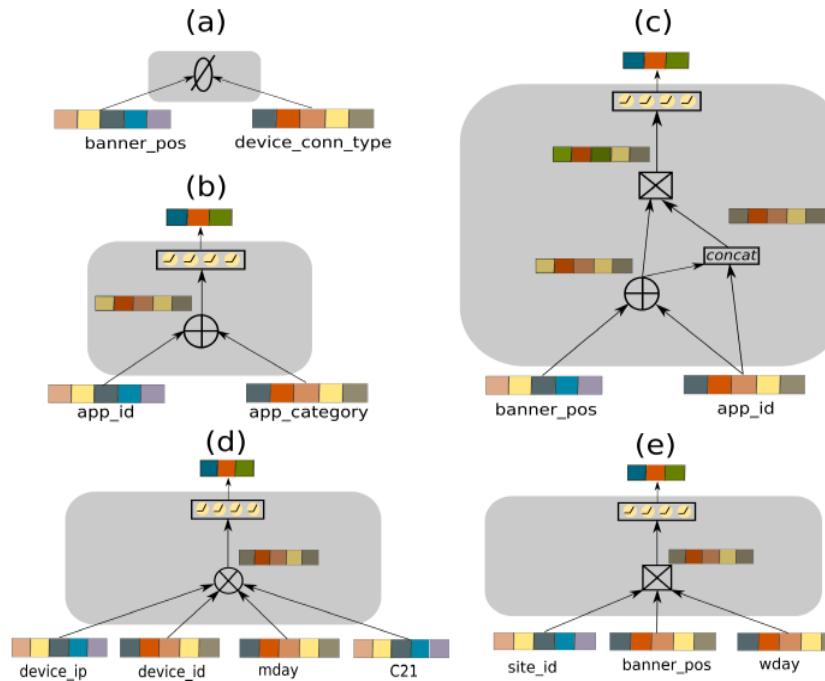


- Vector-level: simple linear algebra operations
- Elementwise: shared nonlinear transformation

# Interaction Function Search-AutoFeature



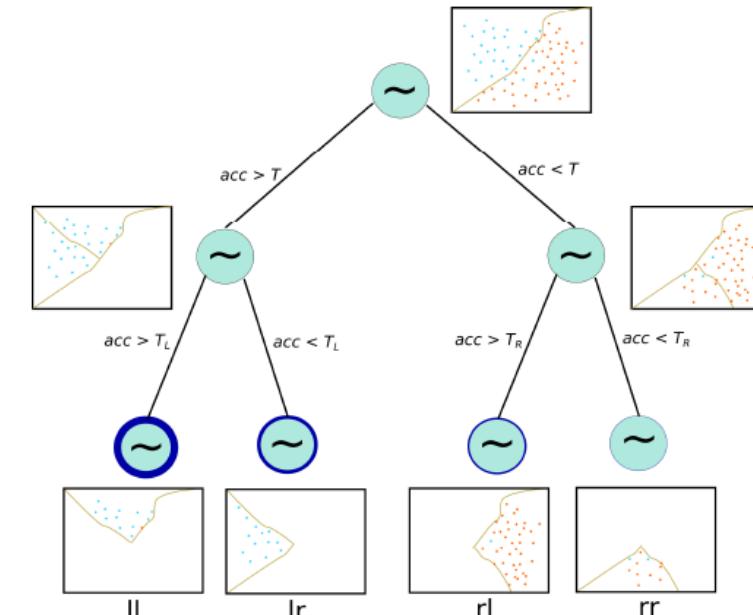
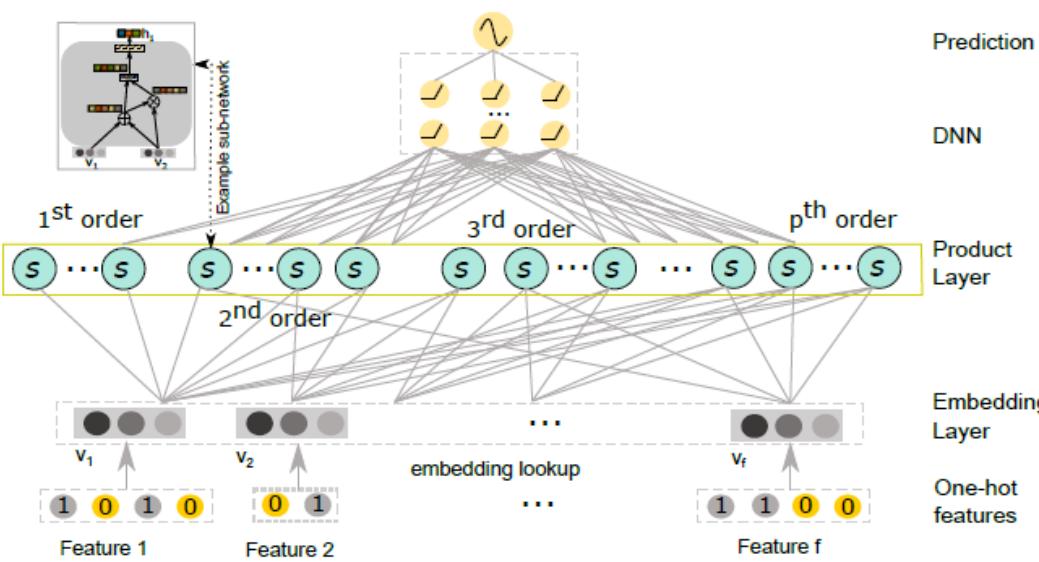
- Not all the feature interactions between each pair of fields need to be modeled.
- Not all the useful feature interactions can be modeled by the same interaction functions



# Interaction Function Search-AutoFeature



- AutoFeature automatically designs a different sub-net for each pair of fields.
- Train a Naïve Bayes Tree to classify different network structures, where the tree tends to classify well-performed network structures into several leaves, such that the next generation can be more effective.

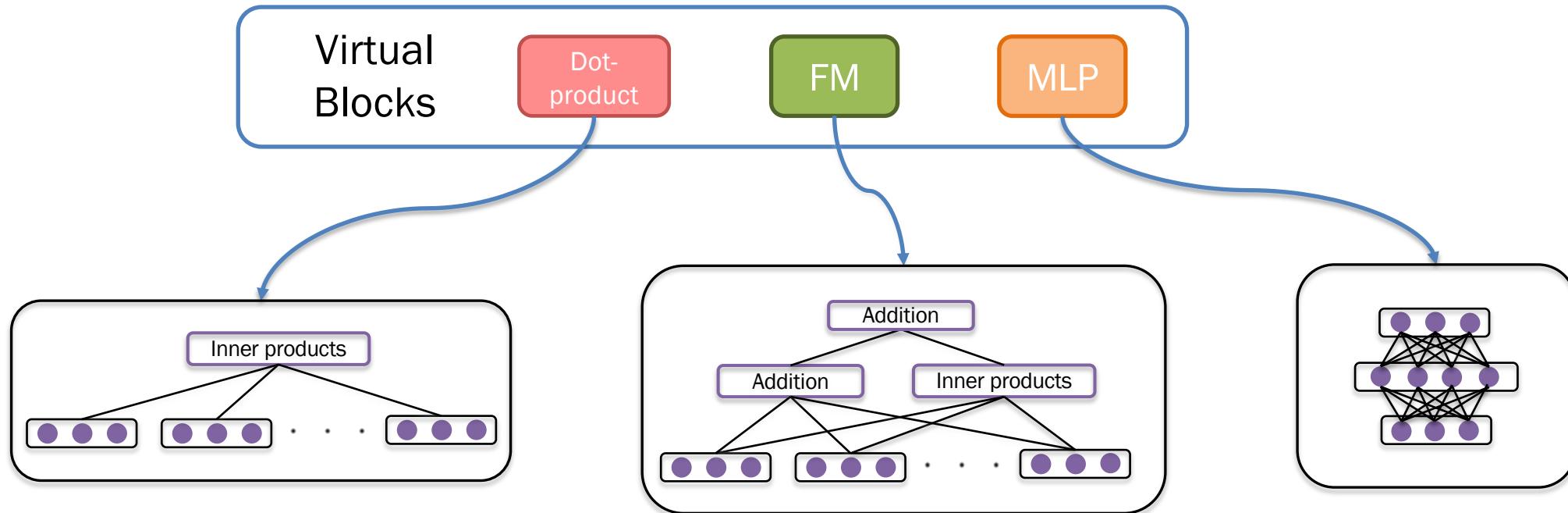


# Interaction Block Search-AutoCTR



## Hierarchical Search Space

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

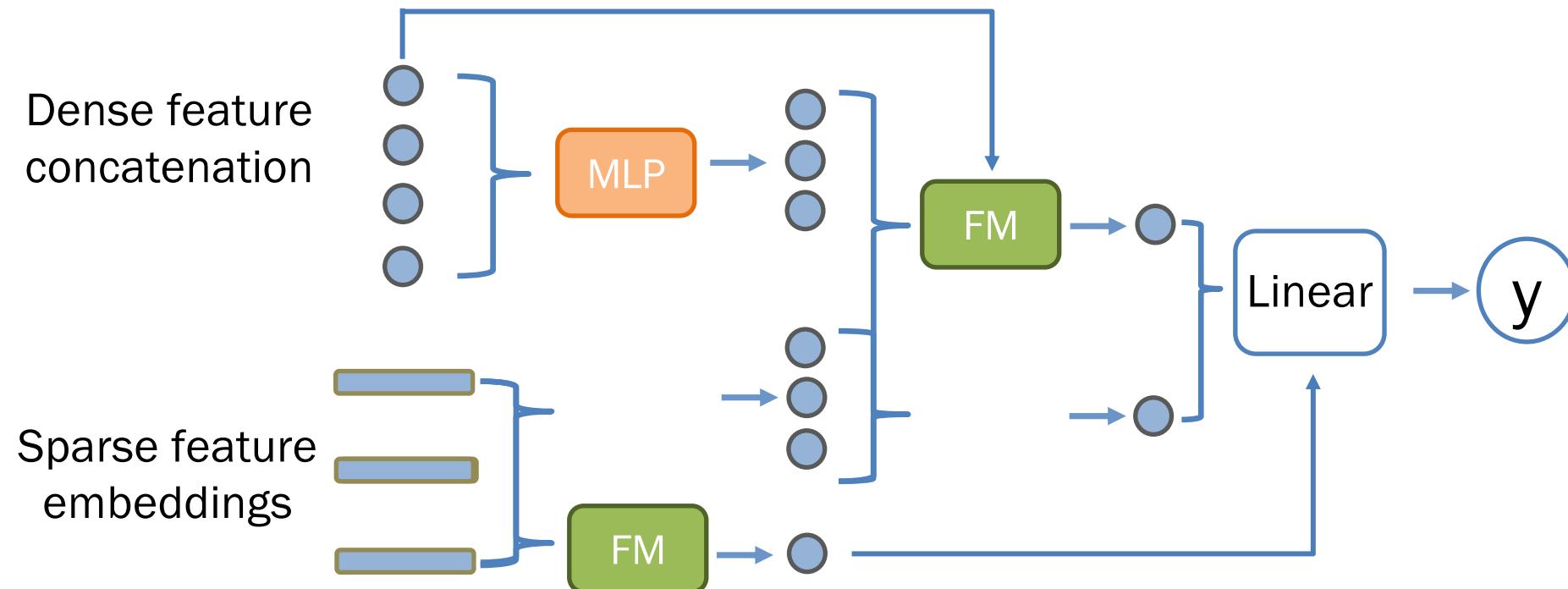


# Interaction Block Search-AutoCTR



## Search space construction

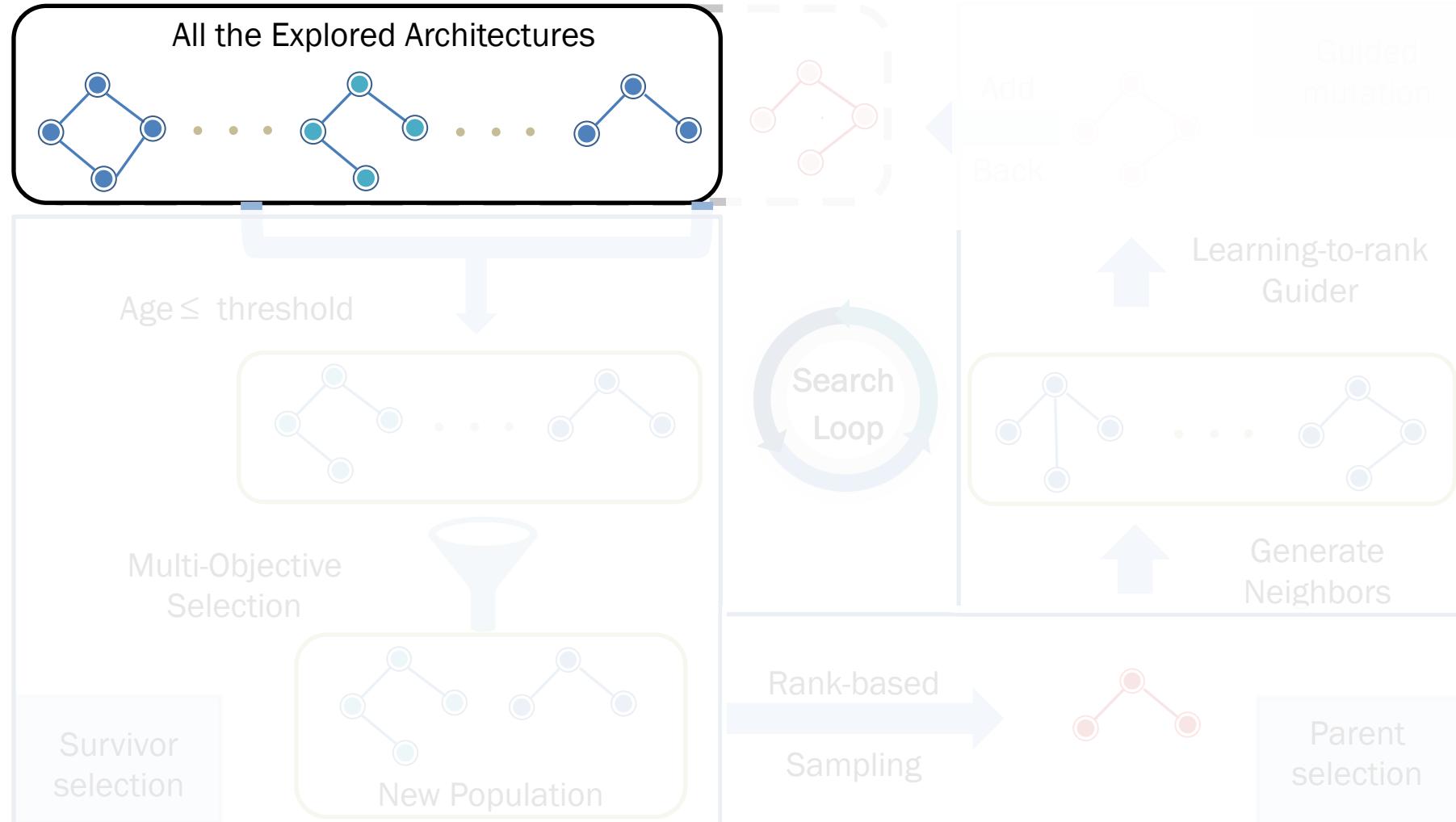
- DAG of virtual blocks and grouped feature embeddings
- Both block hyper-parameters and connection among blocks are to be searched



# Interaction Block Search-AutoCTR



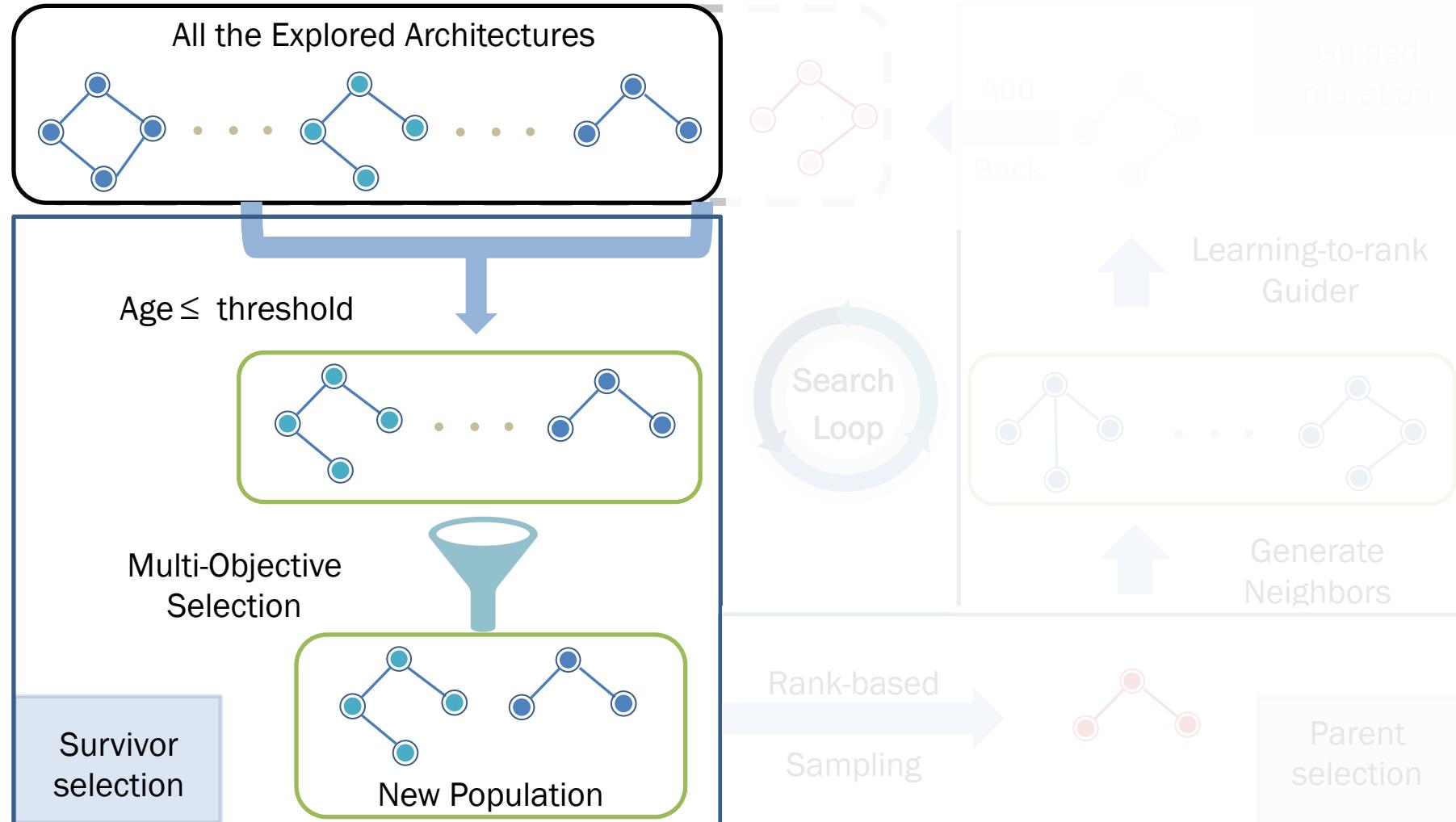
## Multi-Objective Evolutionary Search Algorithm



# Interaction Block Search-AutoCTR



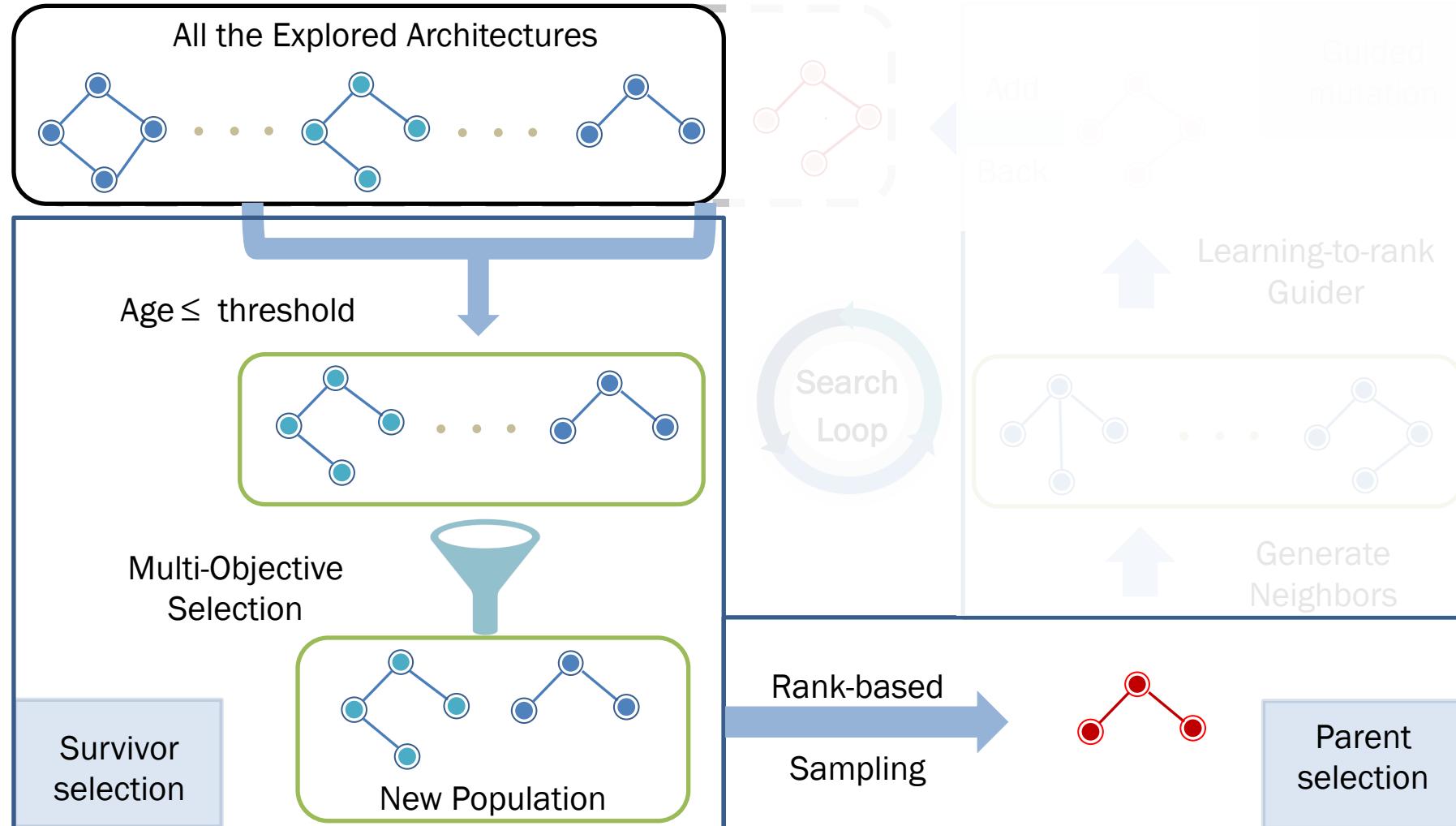
## Multi-Objective Evolutionary Search Algorithm



# Interaction Block Search-AutoCTR



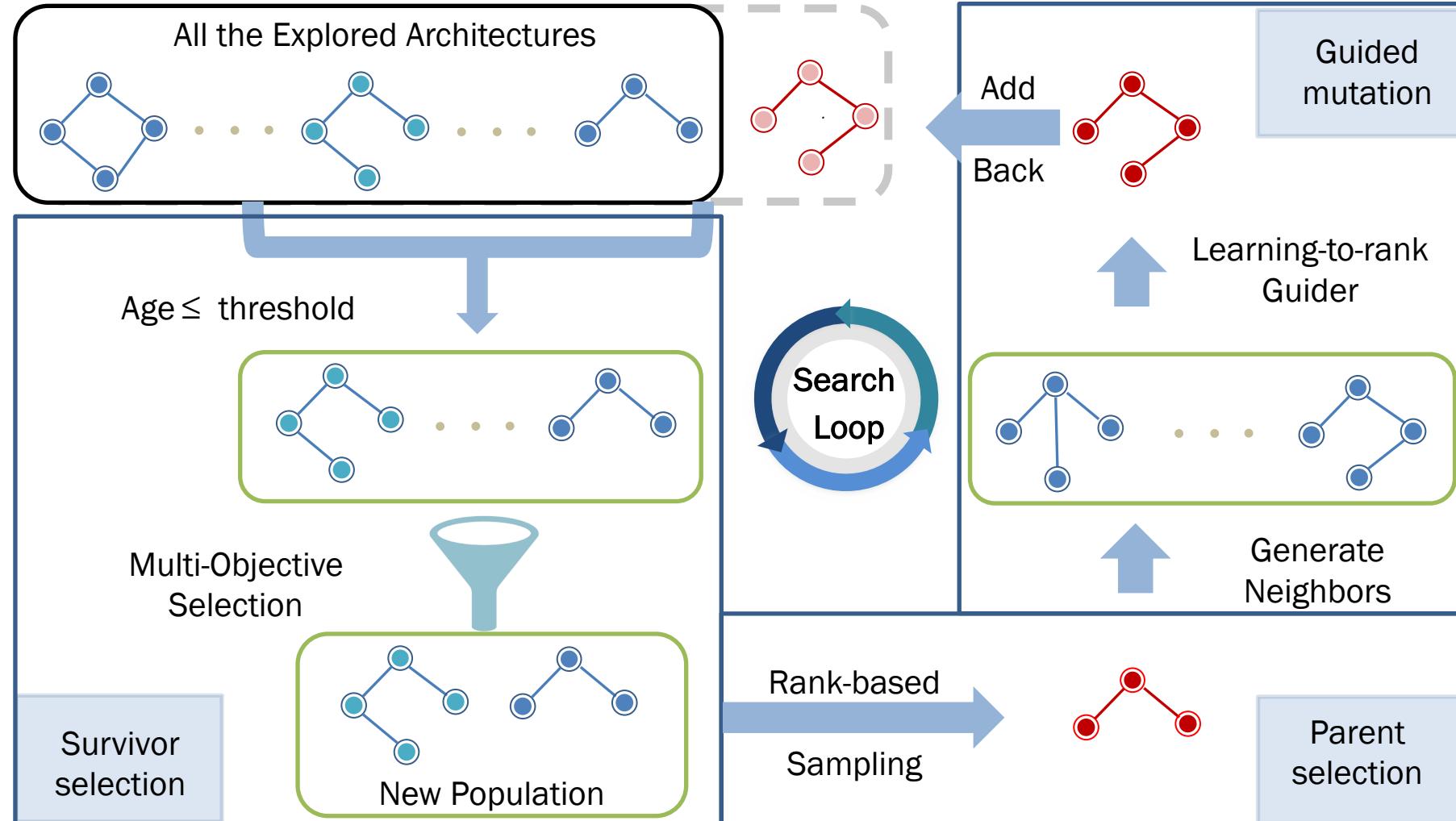
## Multi-Objective Evolutionary Search Algorithm



# Interaction Block Search-AutoCTR



## Multi-Objective Evolutionary Search Algorithm

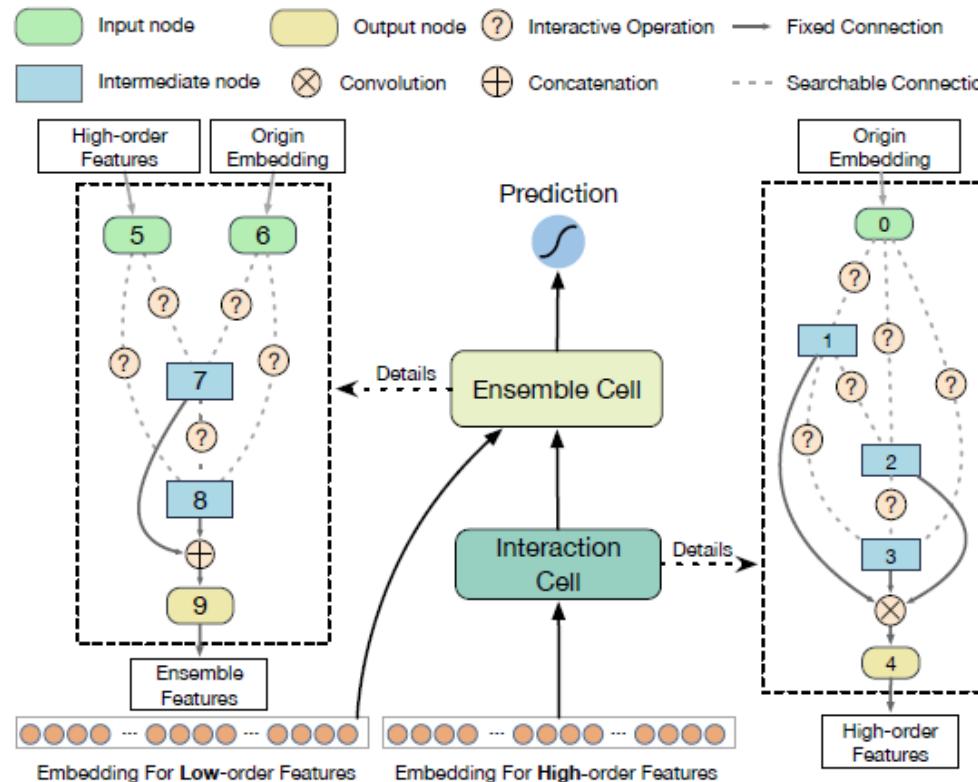


# Interaction Block Search-AutoPI



## Search Space

- The interaction cell formulates the higher-order feature interactions
- The ensemble cell formulates the ensemble of lower-order and higher-order interactions



- Skip-connection<sup>[1]</sup>

$$\mathbf{E}' = o_{\text{skip}}(\mathbf{E}) = \mathbf{E} \in \mathbb{R}^{m \times k}$$

$$\mathbf{E}' = [a'_1 \cdot \mathbf{x}_1, a'_2 \cdot \mathbf{x}_2, \dots, a'_m \cdot \mathbf{x}_m] \in \mathbb{R}^{m \times k}$$

- Self-attention<sup>[3]</sup>

$$\mathbf{E}' = [\mathbf{e}_1^{\text{Res}}, \mathbf{e}_2^{\text{Res}}, \dots, \mathbf{e}_m^{\text{Res}}] \in \mathbb{R}^{m \times k}$$

- FC Layer

$$\mathbf{E}' = \mathbf{e} \cdot \mathbf{W}_{\text{FC}}$$

- FM Layer

$$\begin{aligned} \mathbf{p} &= \{(\mathbf{x}_j \cdot \mathbf{x}_j)\}_{(i,j) \in R_x} \\ \text{s.t. } \mathbf{R}_x &= \{(i,j)\}_{i \in \{1, \dots, m\}, j \in \{1, \dots, m\}, i < j} \end{aligned}$$

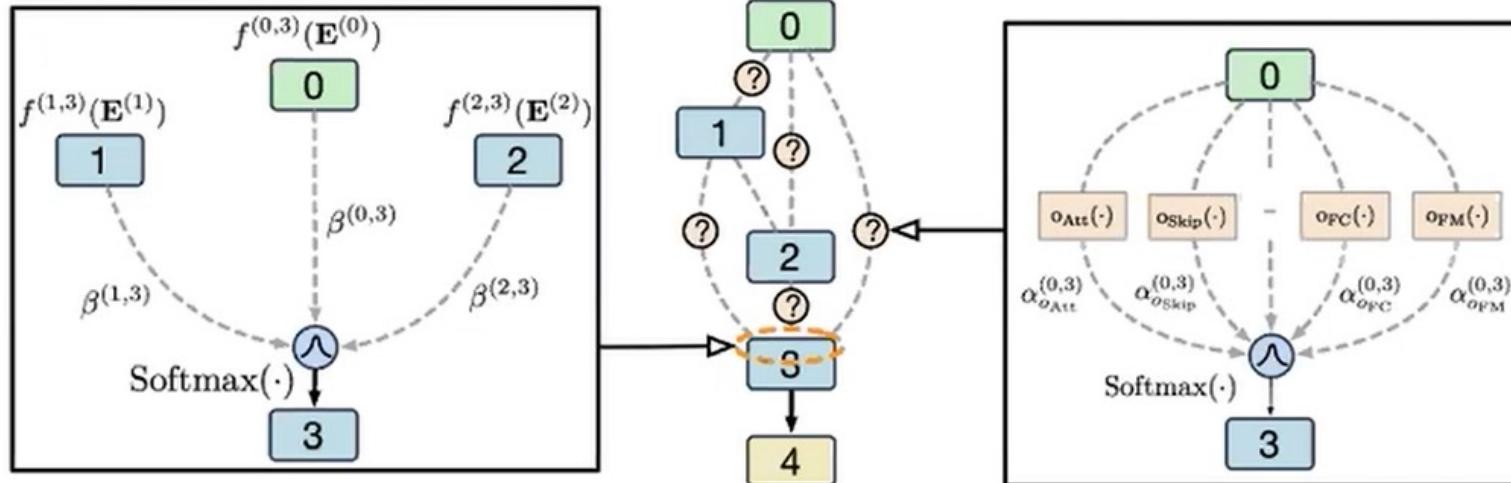
- 1D Conv

$$\left\{ \mathbf{C}^{(i)} \in \mathbb{R}^{1 \times 1 \times m} \right\}_{i \in \{1, \dots, m\}}$$

# Interaction Block Search-AutoPI

## Search Strategy

- Continuous relaxation



Continuous relaxation visualization

By introducing the **operator-level** and **edge-level** architecture parameters for continuous relaxation, a differentiable objective function will be obtained:

$$\min_{\alpha, \beta} \mathcal{L}_{\text{val}}(w^*(\alpha, \beta), \alpha, \beta)$$

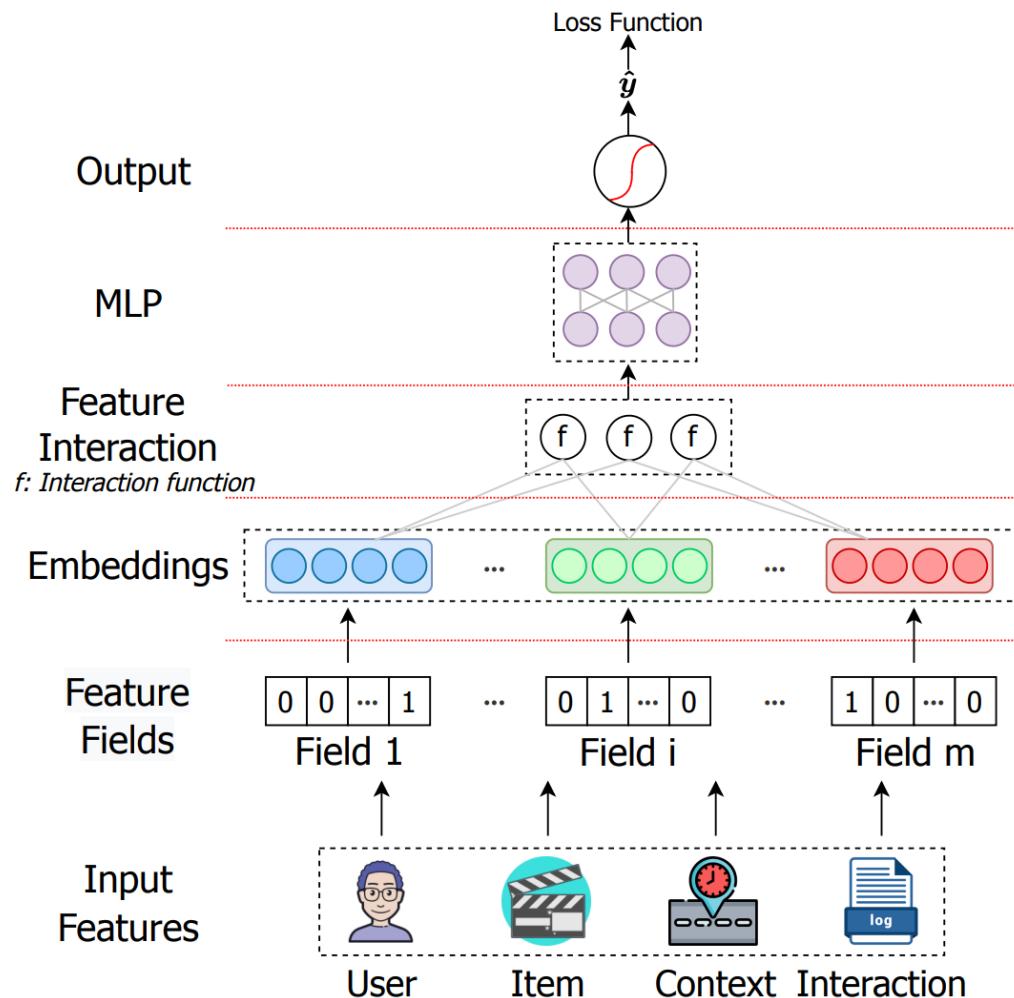
$$\text{s.t. } w^*(\alpha, \beta) = \operatorname{argmin}_w \mathcal{L}_{\text{train}}(w, \alpha, \beta)$$



# Table of Contents

- Introduction
  - Background: Deep Recommender Systems
- Preliminary of AutoML
- DRS Embedding Components
  - Single Embedding Search
  - Group Embedding Search
- DRS Interaction Components
  - Feature Interaction Search
  - Interaction Function Search
  - Interaction Block Search
- **DRS Comprehensive Search & System**
- Conclusion & Future Direction

# Background



- **Comprehensive Search:**
  - Search for several parts of DRS
  - E.g. Embedding size & Feature Interaction function
- **System Design:**
  - Search for architectures other than aforementioned parts
  - E.g. Loss function

# Comprehensive Search & System



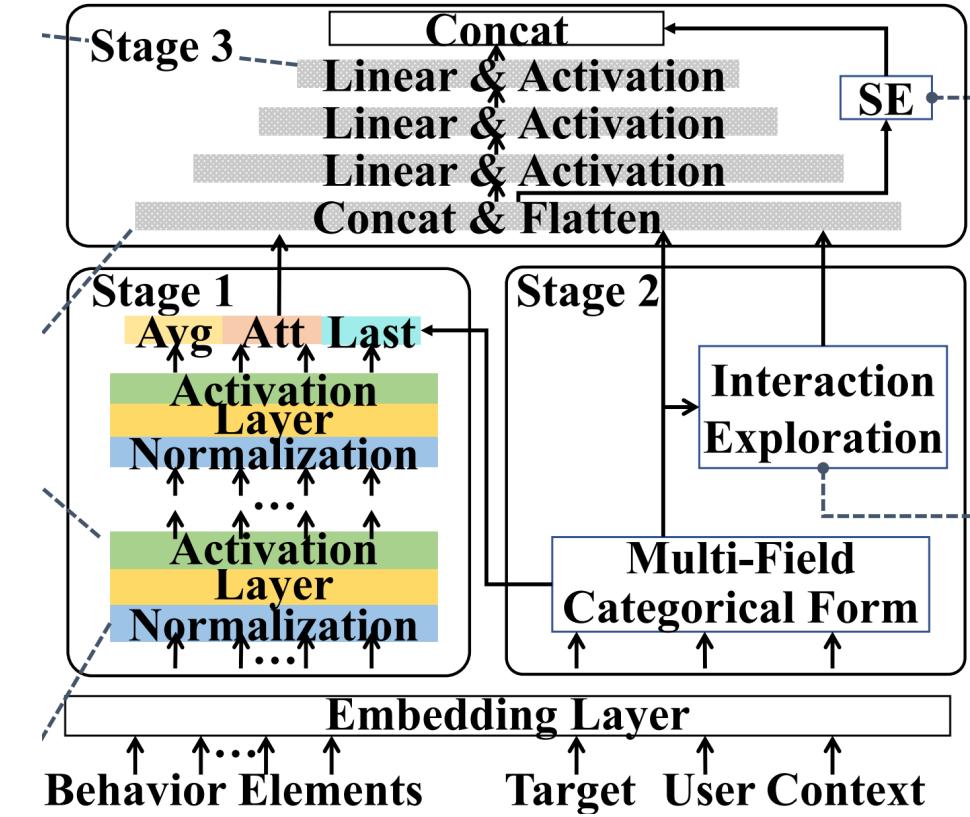
	Model Name	Search Strategy	Search Space
Comprehensive Search	<b>AIM</b>	Gradient	Embedding Dimension, Interaction Function, Feature Interaction
	AutoIAS	Reinforcement Learning	Embedding Dimension, Projection Dimension, Interaction Function, Feature Interaction, MLP Layer & Dimension
	<b>AMEIR</b>	One-shot Random Search	Sequential model, Feature interaction, MLP Dimension
System Design	<b>AutoLoss</b>	Gradient	Optimization: Loss Function
	AutoGSR	Gradient	Structure Design: GNN Architecture
	AutoFT	Gradient	Parameter Tuning: Fine-Tune or Not (For pre-trained models)

- Motivation:
  - RS could be divided into 3 parts: Behavior modeling, feature interaction, MLP
  - Hard to find a unified model for all scenarios
  - Existing NAS methods for Recommendation are restricted
- Problem:
  - Automatically find the **complete** recommendation models for adapting **various** recommender tasks
- Solution:
  - Three-step one-shot random search, progressively construct the architecture

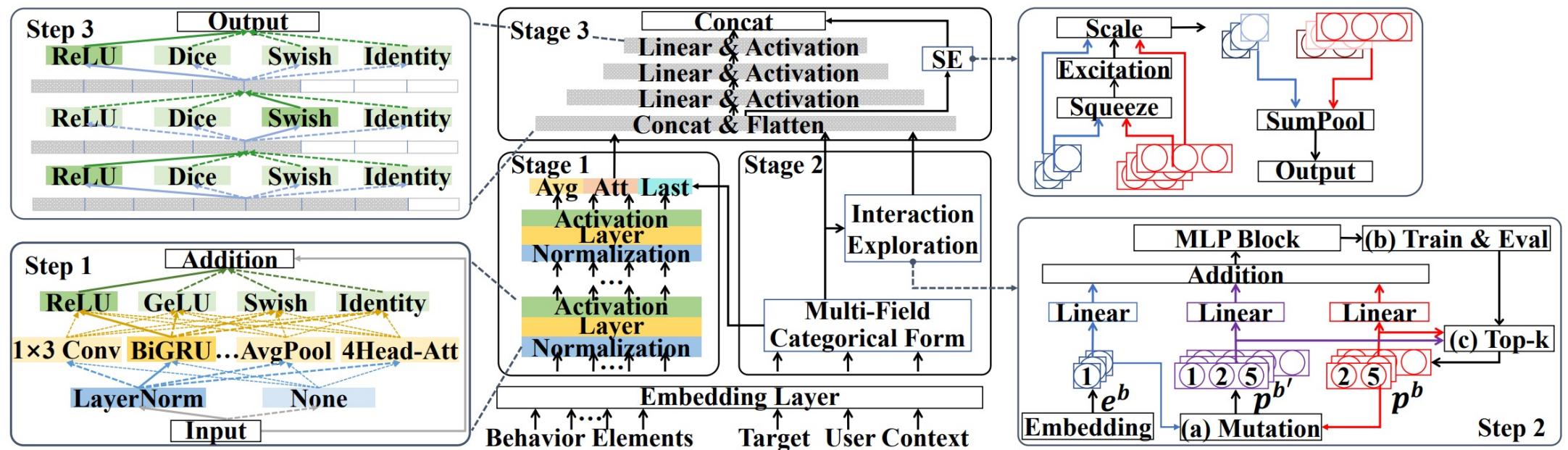
# AMEIR – Search Space

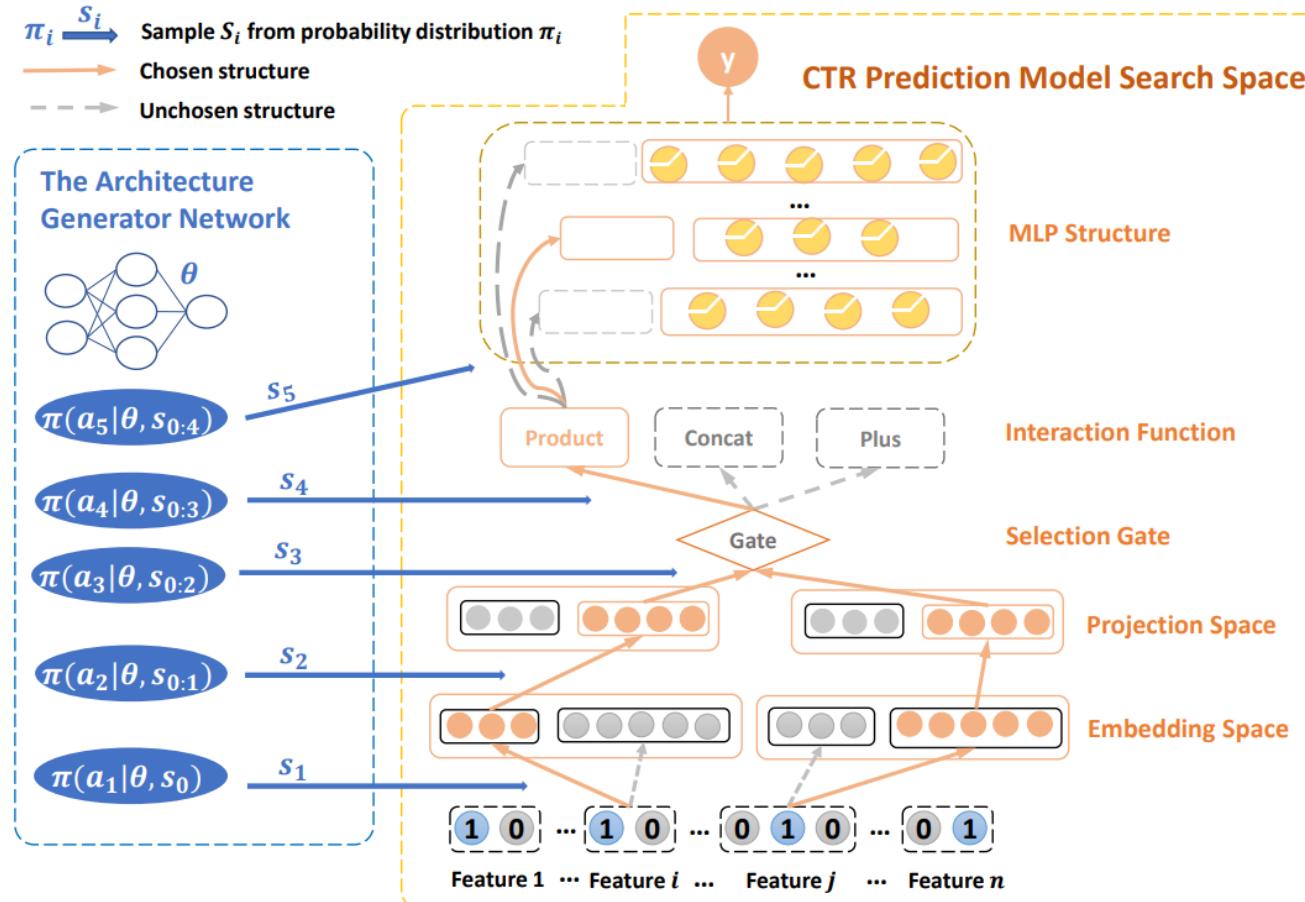


- Stage 1: (Behavior Modeling for Sequential features)
  - Search for a fixed number of layers ( $L$ )
  - Normalization: {Layer Norm, None}
  - Layer: {Conv, Recur, Pooling, Attention} & Zero
  - Activation: {ReLU, GeLU, Swish, Identity}
  - Output: (Average, Attn, Last) please refer to the article
- Stage 2: (Interaction Exploration for non-Sequential features)
  - Interaction function: Hadamard Product (Fixed)
  - Only Search for feature combination for interaction
- Stage 3: (MLP investigation)
  - MLP Dimension:  $\{0.1, 0.2, \dots, 1.0\}$  of input size
  - Activation: {ReLU, Swish, Identity, Dice}



- Overall Search Strategy: One-shot random search
- Step 1: Using a predefined MLP, search for the optimal architecture.
- Step 2: Combined with SMBO, progressively expand the interaction sets, also use a predefined MLP.
- Step 3: Design a weight matrix of maximal dimension to realize one-shot search



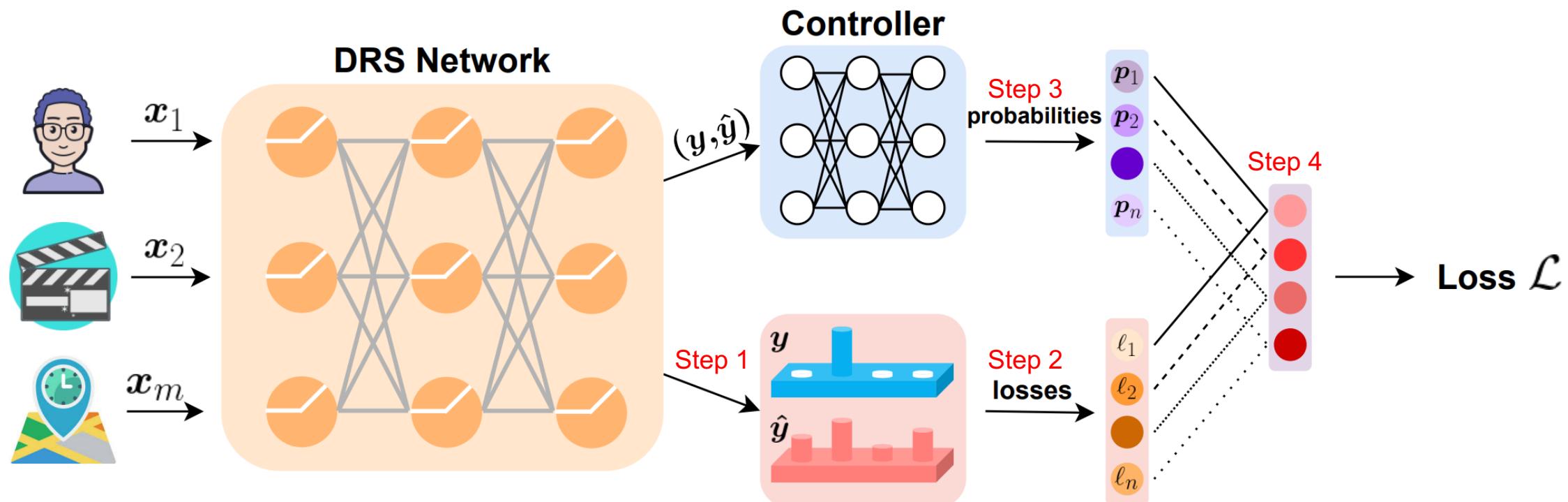


- Search Space:
  - Embedding: Embedding Size
  - Projection: Unified Size for interaction
  - Feature selection: Pairs of feature
  - Interaction Function: {Product, ...}
  - MLP:
    - The number of layers
    - Layer dimensions
- Strategy: Reinforcement Learning

- Motivation:
  - Traditional Setting: Predefined and fixed loss function, sub-optimal
  - Exhaustively or manually searched fused loss: Costly, neglecting the difference between data examples
- Problem:
  - automatically and adaptively assign the appropriate loss function for each data example
- Solution:
  - Design a controller network to adaptively adjust the probabilities over multiple loss functions for different data examples

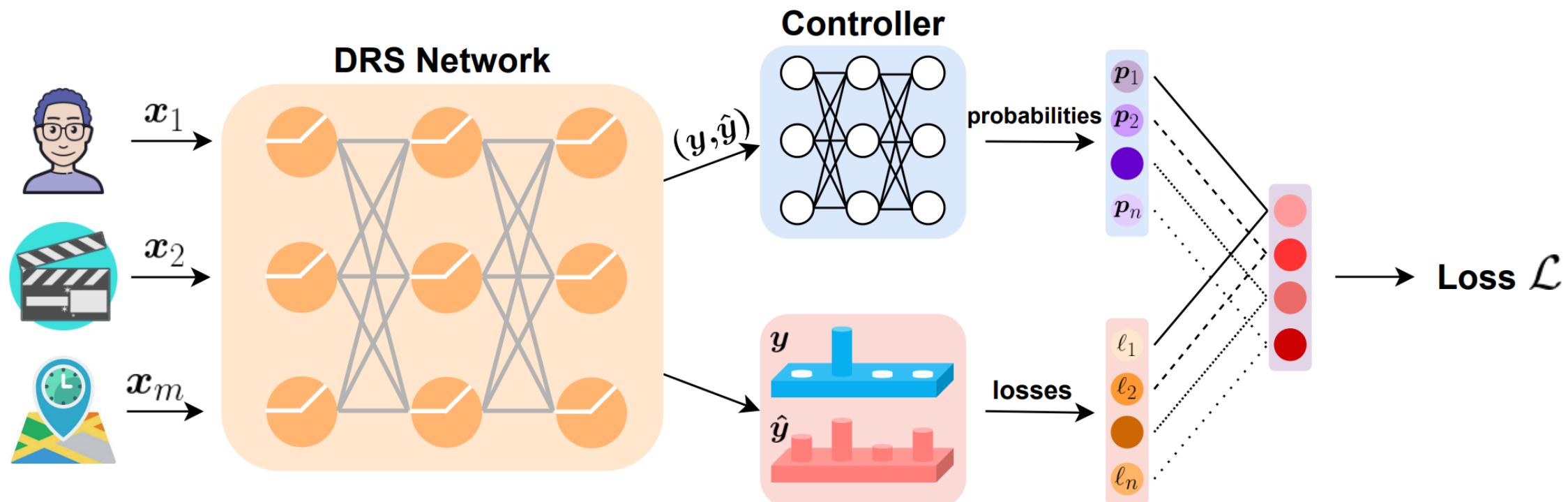
# AutoLoss – Forward-propagation

- Step 1: The DRS makes predictions
- Step 2: Calculate candidate losses
- Step 3: The controller generate weights(probabilities) according to predictions.
- Step 4: Calculate the overall Loss (Weighted sum)

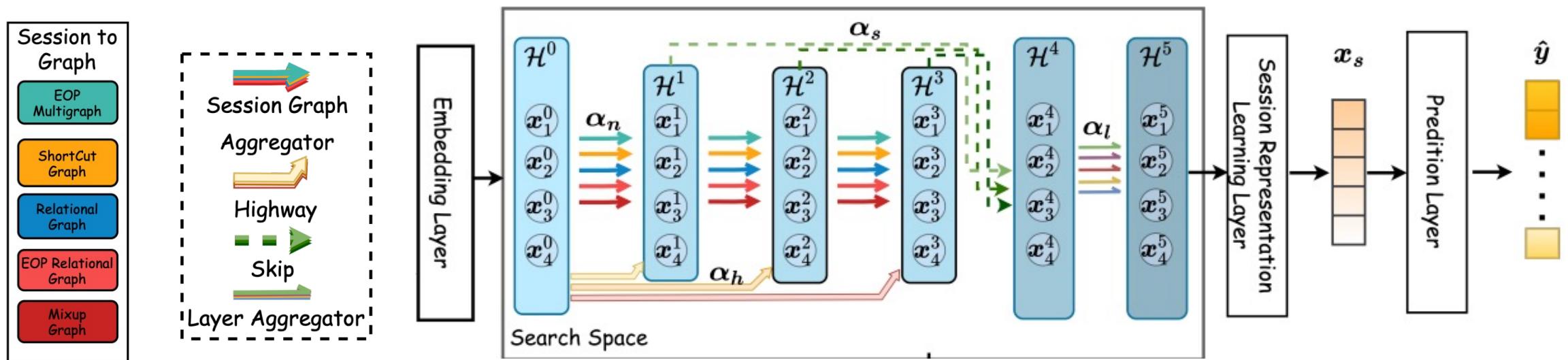


# AutoLoss – Backward-propagation

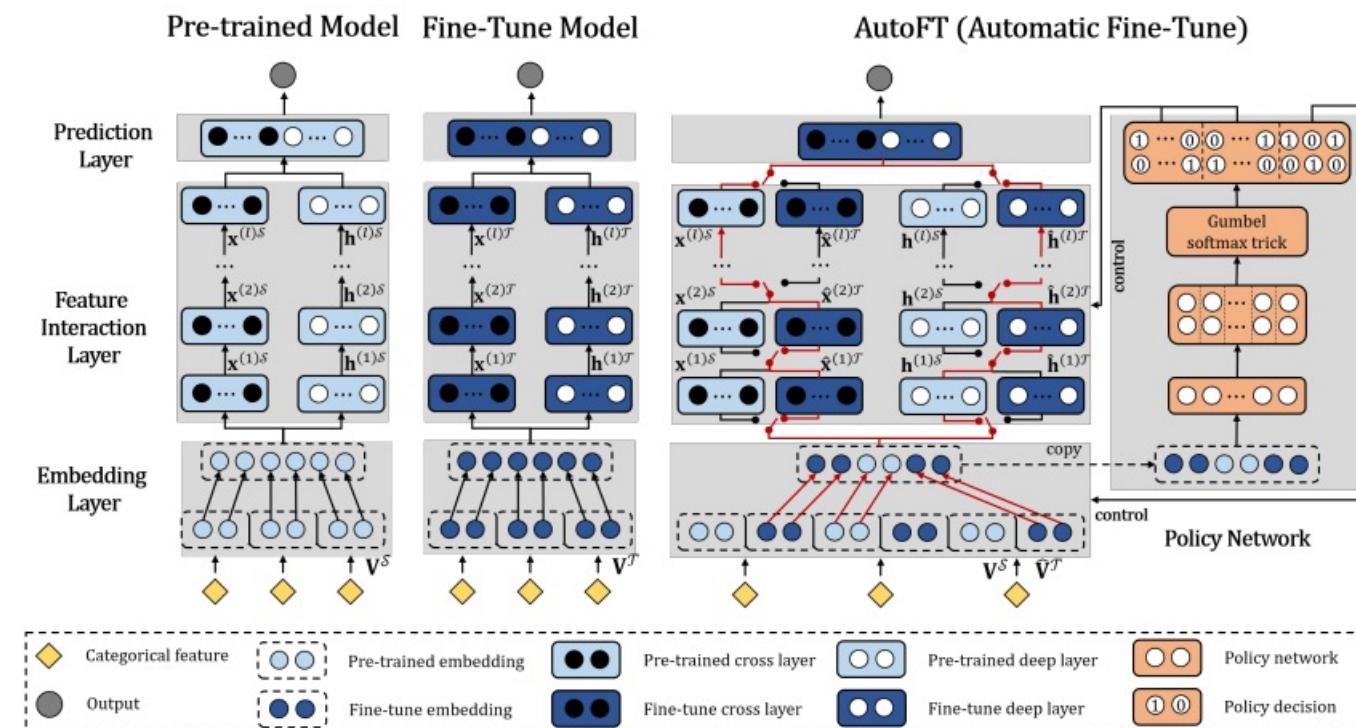
- Parameters of the DRS network: Updated based on **Training Error**
- Parameters of the Controller: Updated based on **Validation Error**



- Target: Search GNN architectures for Session-based Recommendation
- Methodology:
  - Search Space:
    - Session Aggregation: 5 popular graph types.
    - Layer Aggregation: mean, max, concat, sum & highway & skip.
  - Strategy: Continuous relaxion & Gradient



- Target: Automatically which parts of pre-trained model should be fine-tuned
- Search Space:
  - Field-wise transfer: use fine-tuned embedding or not
  - Layer-wise transfer: Whether to fine-tune cross & deep layers
- Strategy: Continuous relaxion & Gradient





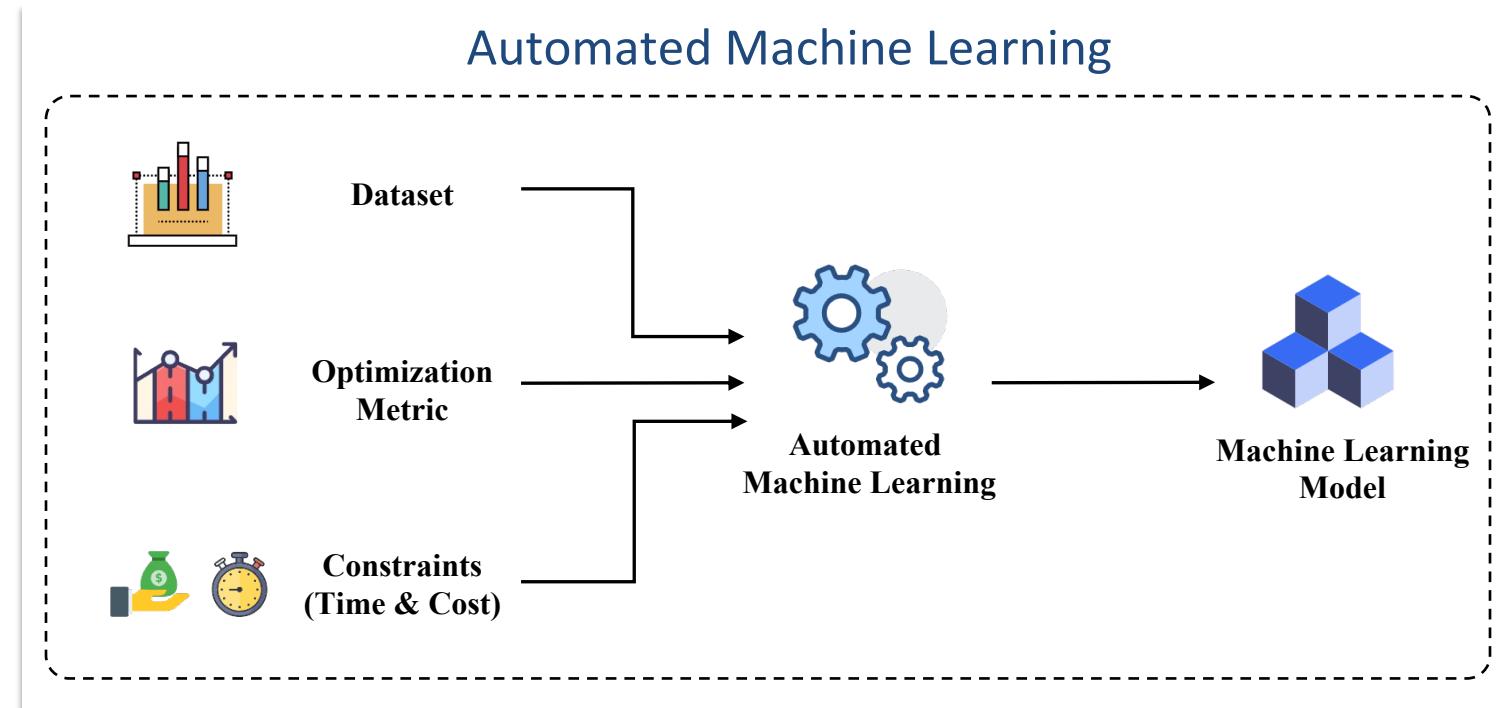
# Table of Contents

- Introduction
  - Background: Deep Recommender Systems
- Preliminary of AutoML
- DRS Embedding Components
  - Single Embedding Search
  - Group Embedding Search
- DRS Interaction Components
  - Feature Interaction Search
  - Interaction Function Search
  - Interaction Block Search
- DRS Comprehensive Search & System
- Conclusion & Future Direction

# Conclusion

Automated Machine Learning contribute to improving the performance of recommender systems in a data-driven manner.

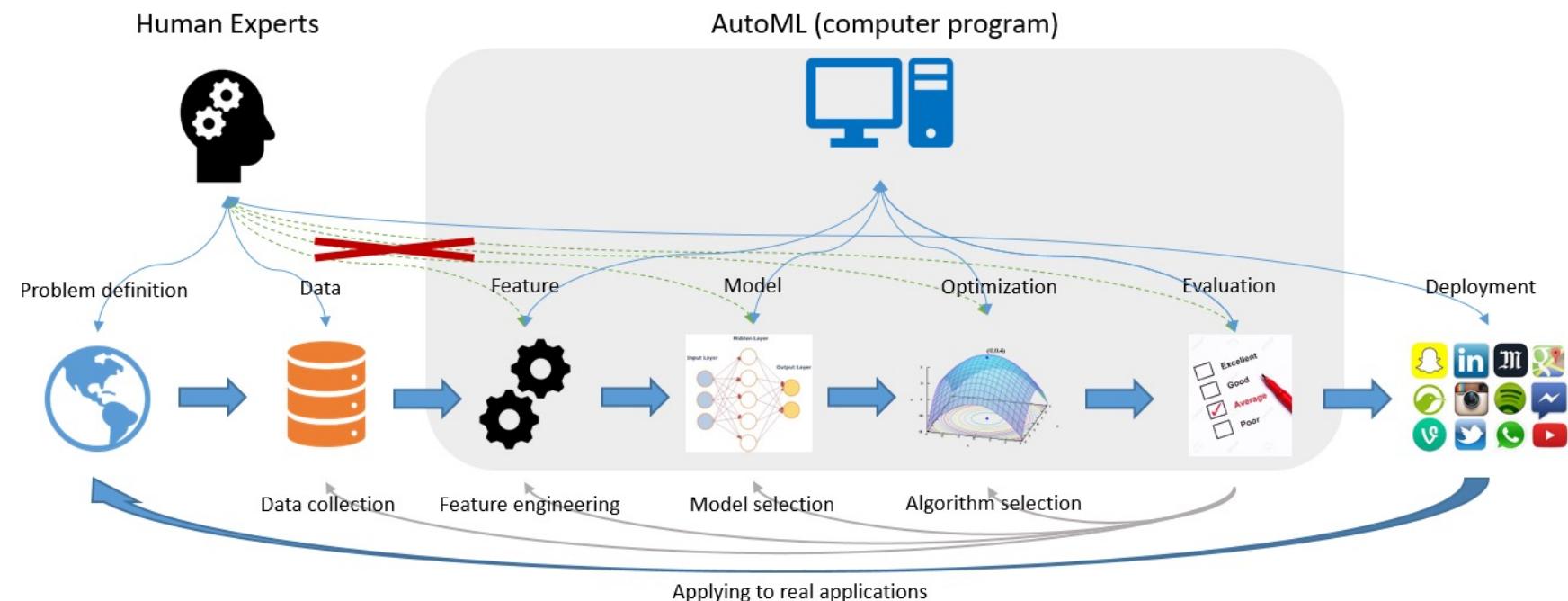
- Automatically search embedding dimensions to better model **feature representations**
- Automatically design deep networks to better capture **feature interactions**
- Automatically design **comprehensive system architectures** to better improve performance



# Conclusion

AutoML advantages:

- Less expert knowledge
- Saving time and efforts
- Different data → different architectures



## 1) Feature Selection

- Recommendation scenario contains a plenty of features (user, item, context, and cross features)
- Automatically search optimal feature sub-set under hardware and latency constraints

## 2) Model Selection

- Search different sub-models/sub-architectures according to different requests adaptively

## 3) User Behavior Modeling

- User history behaviors contain different dimensions of interests
- Automatically retrieve beneficial history behaviors for modeling user preference

## 4) Multi-task learning

- Multi-task learning is one of the most important techniques in industry recommendation for considering different revenue targets (e.g., ctr, cvr, vv)
- Designing an automatic algorithm for recommendation based on multi-task learning

## 5) GNNs-based Recommendations

- Exploring the combination of AutoML and GNNs provides great opportunities to further boost the performance of GNNs-based recommendation methods.

## 6) Multi-Modality Recommendations

- Designing an optimal algorithm via automated machine learning techniques to advance deep multimodal learning for recommendation.



# Thank you!