

Interactions Expressed in Sequences and Graphs

Sequential Recommendations and GNN-based RS

Klaudia Balcer

Advanced Data Mining



Contents

1 Sequential Recommendations

- GRU4Rec
- SASRec

2 Graph-based RS

- LightGCN
- [uncovered] SR-GNN, TAGNN

3 [uncovered] Augmentations for Interactions

- [uncovered] SimGCL
- [uncovered] DiffuASR

4 [uncovered] Matrix Form – revisited

- [uncovered] DiffRec

Sequential Recommendations

Sequential Recommendations

GRU4Rec

GRU

much to update the hidden state of the unit. The activation of the GRU is a linear interpolation between the previous activation and the candidate activation \hat{h}_t :

$$h_t = (1 - z_t)h_{t-1} + z_t \hat{h}_t \quad (2)$$

where the update gate is given by:

$$z_t = \sigma(W_z x_t + U_z h_{t-1}) \quad (3)$$

while the candidate activation function \hat{h}_t is computed in a similar manner:

$$\hat{h}_t = \tanh(W x_t + U(r_t \odot h_{t-1})) \quad (4)$$

and finally the reset gate r_t is given by:

$$r_t = \sigma(W_r x_t + U_r h_{t-1}) \quad (5)$$

GRU4Rec

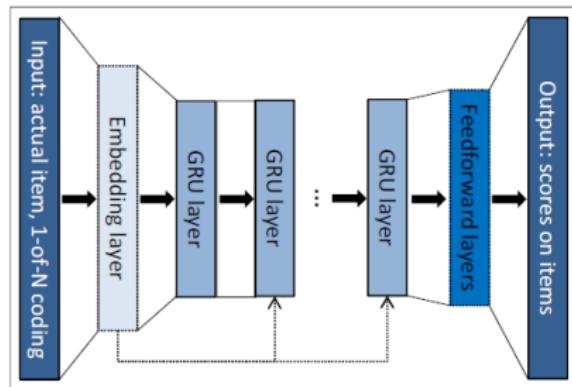


Figure 1: General architecture of the network. Processing of one event of the event stream at once.

- session-based recommendations (when matrix factorization is not suitable)

<https://arxiv.org/abs/1511.06939>

GRU4Rec

Session1 $i_{1,1} | i_{1,2} | i_{1,3} | i_{1,4}$

Session2 $i_{2,1} | i_{2,2} | i_{2,3}$

Session3 $i_{3,1} | i_{3,2} | i_{3,3} | i_{3,4} | i_{3,5} | i_{3,6}$

Session4 $i_{4,1} | i_{4,2}$

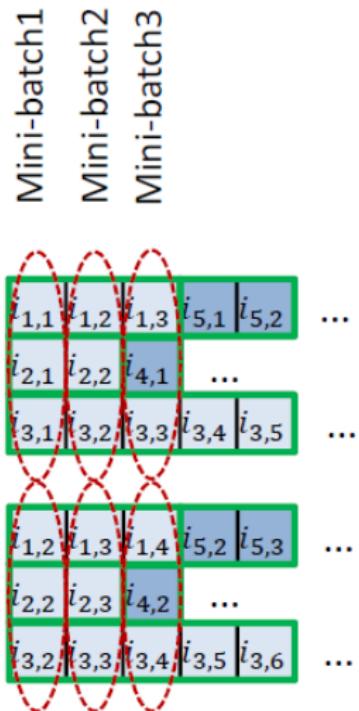
Session5 $i_{5,1} | i_{5,2} | i_{5,3}$

...



Input

Output

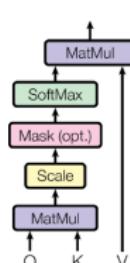


Sequential Recommendations

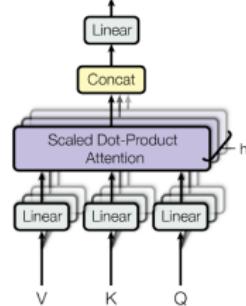
SASRec

Transformer

Scaled Dot-Product Attention



Multi-Head Attention



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Q, K V - Queries, Keys, and Values matrices
 order in sequence → positional embeddings

https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fb053c1c4a845aa-Abstract.html

Self-Attentive Sequential Recommendation

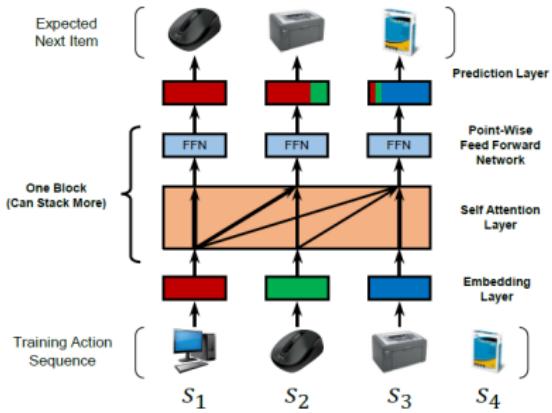


Figure 1: A simplified diagram showing the training process of SASRec. At each time step, the model considers all previous items, and uses attention to 'focus on' items relevant to the next action.

*The goal of our work is to balance [...], by proposing a sequential model based on self-attention (SASRec) that allows us to capture **long-term semantics** (like an RNN), but, using an attention mechanism, makes its **predictions based on relatively few actions** (like an MC).*

SASRec

- Embeddings Layer (fixed length: cutting/padding):

$$\hat{V}_s = \begin{bmatrix} V_{s_1} + P_1 \\ V_{s_2} + P_2 \\ \vdots \\ V_{s_n} + P_n \end{bmatrix}$$

Positional Embeddings P are trainable.

- Attention block:

- Self-Attention Layer:

$$S := SA(\hat{V}_s) = \text{Attention}(Q : \hat{V}_s W^Q, K : \hat{V}_s W^K, V : \hat{V}_s W^V)$$

Causality: modify the attention by forbidding all links between Q_i and K_j if $j > i$.

- Feed forward network: linear -> ReLU -> linear on S
- To prevent overfitting authors use: layer normalization, dropout and residual connections.

- The predictions are made on the basis of the dot product.
- Binary cross-entropy loss (1 random negative sample).
- Attention weights give us some interpretability (see the paper).

Graph-based RS

Graph-based RS

LightGCN

Vanilla Gradient Neural Network (recap)

- $\mathcal{G} = (\mathcal{V}, \mathcal{E}, X^{\mathcal{V}}, X^{\mathcal{E}})$ - a graph is a tuple of vertices, edges, vertices features and edges features.
- The target is to provide meaningful node embedding matrix V .
- In each layer we have two major phases:

- ① Aggregation:

$$\mathbf{m}_v^{(l)} = f \left(X_v^{\mathcal{V}}, X_{\mathcal{N}_v}^{\mathcal{E}}, X_{\mathcal{N}_v}^{\mathcal{V}}, V_{\mathcal{N}_v}^{(l-1)} \right)$$

- ② Message passing:

$$\mathbf{v}_v^{(l)} = g \left(\mathbf{v}_v^{(l-1)}, \mathbf{m}_v^{(l)} \right)$$

https://cs.mcgill.ca/~wlh/comp766/files/chapter4_draft_mar29.pdf

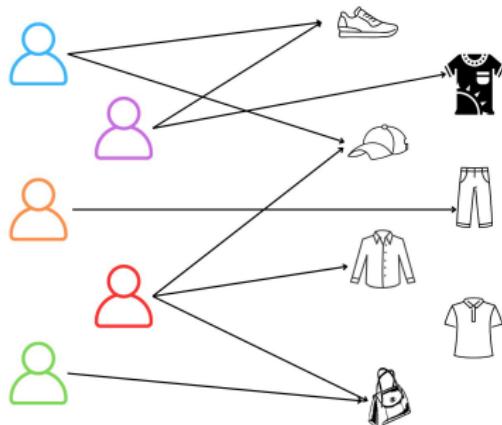
Graph Convolutional Network (recap)

$$\mathbf{V}^{(l)} = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} \mathbf{V}^{(l-1)} \mathbf{W}^{(l)} \right)$$

- $\tilde{A} = A + I_N$: the adjacency matrix of the undirected graph \mathcal{G} with added self-connections (I_N – the identity matrix).
- $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$: the diagonal degree matrix of \tilde{A} .
- $\mathbf{W}^{(l)}$: a layer-specific trainable weight matrix.
- $\sigma(\cdot)$: an activation function, such as the ReLU function, $\text{ReLU}(\cdot) = \max(0, \cdot)$.
- $\mathbf{V}^{(l)} \in \mathbb{R}^{N \times D}$: the matrix of activations in the l^{th} layer.
- $\mathbf{V}^{(0)} = X^{\mathcal{V}}$: the input feature matrix.

$$\mathbf{V}_v^{(l)} = \sigma \left(\frac{1}{\sqrt{|\mathcal{N}_v| + 1}} \sum_{w \in \mathcal{N}_v \cup \{v\}} \frac{1}{\sqrt{|\mathcal{N}_w| + 1}} \mathbf{V}_w^{(l-1)} \mathbf{W}^{(l)} \right)$$

Data Structure



- $\mathcal{V} = (\mathcal{U}, \mathcal{I})$ – we have nodes for both users and items
- $\mathcal{E} = \{(u, i) : \text{user } u \text{ interacted with item } i\}$
- $A = \begin{bmatrix} 0 & R \\ R^\top & 0 \end{bmatrix}$ – adjacency matrix, R – binary interaction matrix
- We need both item and feature embeddings. We use the standard procedure with dot-product scores and BPR loss.

LightGCN - Key Idea

- ~~node and edge features~~
- ~~non-linear activations~~
- neighbourhood aggregation

$$\mathbf{V}^{(l)} = \left(D^{-\frac{1}{2}} A D^{-\frac{1}{2}} \right) \mathbf{V}^{(l-1)}$$

- self-connection obtained after going through all layers:

$$V = \frac{1}{L+1} \sum_{l=0}^L V^{(l)}$$

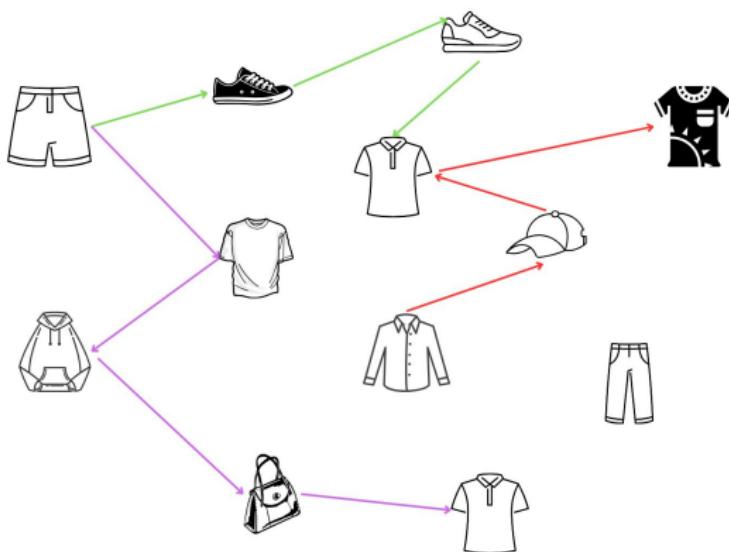
- the only trainable parameters are embedding vectors

<https://arxiv.org/abs/2002.02126> In the paper you can also find:

- > an ablation study on previous Neural Graph Collaborative Filtering model
- > comparison to similar architectures
- > experiments

Graph-based RS
[uncovered] SR-GNN, TAGNN

Data Structure



- $\mathcal{V} = \mathcal{I}$ – we consider only items as nodes as the users are anonymous (and cannot *come back*)
- an edge connects two items consecutive in a sequence

SR-GNN model architecture

$$\mathbf{a}_{m,j}^{(g)} = \mathbf{A}_{m,j} [\mathbf{v}_{m,1}^{(g-1)}, \mathbf{v}_{m,2}^{(g-1)}, \dots, \mathbf{v}_{m,L_m-1}^{(g-1)}]^\top \mathbf{H} + \mathbf{b}, \quad (1)$$

$$\mathbf{z}_{m,j}^{(g)} = \sigma \left(\mathbf{W}_z \mathbf{a}_{m,j}^{(g)} + \mathbf{U}_z \mathbf{v}_j^{(g-1)} \right), \quad (2)$$

$$\mathbf{r}_{m,j}^{(g)} = \sigma \left(\mathbf{W}_r \mathbf{a}_{m,j}^{(g)} + \mathbf{U}_r \mathbf{v}_j^{(g-1)} \right), \quad (3)$$

$$\tilde{\mathbf{v}}_j^{(g)} = \tanh \left(\mathbf{W}_o \mathbf{a}_{m,j}^{(g)} + \mathbf{U}_o \left(\mathbf{r}_{m,j}^{(g)} \odot \mathbf{v}_j^{(g-1)} \right) \right), \quad (4)$$

$$\mathbf{v}_j^{(g)} = \left(1 - \mathbf{z}_{m,j}^{(g)} \right) \odot \mathbf{v}_j^{(g-1)} + \mathbf{z}_{m,j}^{(g)} \odot \tilde{\mathbf{v}}_j^{(g)}, \quad (5)$$

where $\mathbf{A}_{m,j} \in \mathbb{R}^{1 \times 2N}$ is the vector of two columns of the adjacency matrix corresponding to the j th row item, $\mathbf{z}_{m,j}^{(g)}$ is the reset gate, $\mathbf{r}_{m,j}^{(g)}$ is the update gate, $[...]$ is concatenation, σ is the sigmoid function, and \odot is the element-wise multiplication.

<https://dl.acm.org/doi/abs/10.1609/aaai.v33i01.3301346>

SR-GNN model architecture

Then, we need to provide a session embedding s_m . It will be a transformation of a global and a local embedding. The local embedding is just the embedding of the last item in the session. The global embedding is an aggregation of all embeddings in the session. The whole procedure can be described as follows:

$$s_m^l = v_{m,L_m-1}, \quad (6)$$

$$\alpha_j = q^T \sigma(W_1 v_{m,L_m-1} + W_2 v_{m,j} + c), \quad (7)$$

$$s_m^g = \sum_{j=1}^{L_m-1} \alpha_j v_{m,j}, \quad (8)$$

$$s_m = W_3[s_m^l, s_m^g], \quad (9)$$

where q, W_1, W_2 control the weights, W_3 reduces the dimensionality, all of which are learnable parameters.

descriptions from my master thesis

[uncovered] Augmentations for Interactions

Augmentations in Recommender Systems

Data Augmentation for Sequential Recommendation

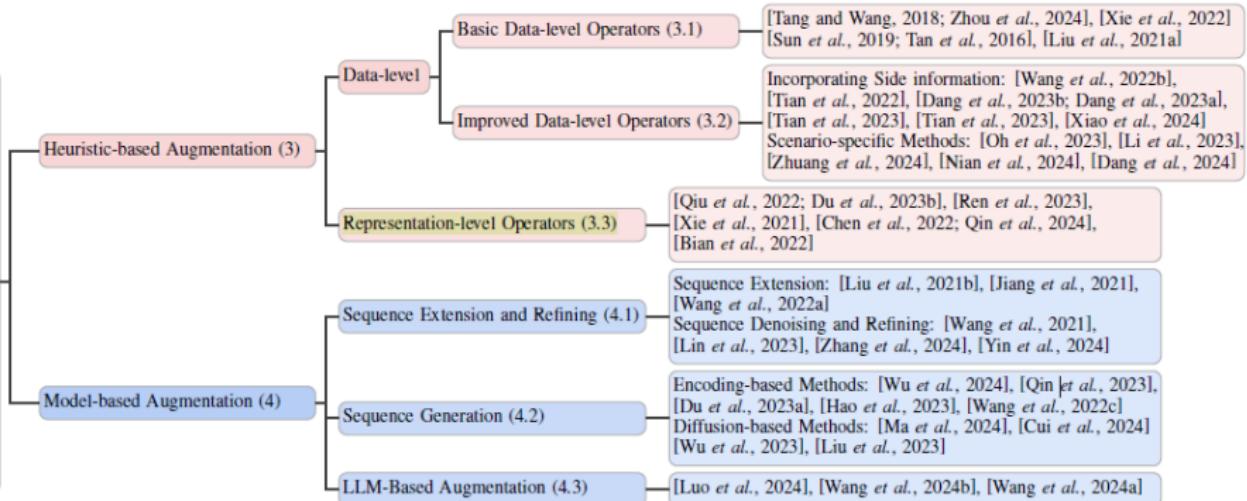


Figure 1: The taxonomy of data augmentation for sequential recommendation.

<https://arxiv.org/abs/2409.13545>

[uncovered] Augmentations for Interactions

[uncovered] SimGCL

SimGCL – noise

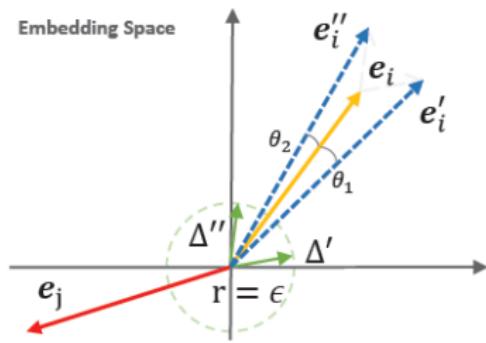


Figure 3: An illustration of the proposed random noise-based data augmentation in \mathbb{R}^2 .

- noise augmentations result in even embedding distribution

<https://dl.acm.org/doi/10.1145/3477495.3531937>

InfoNCE loss function

$$\mathcal{L}_{joint} = \mathcal{L}_{rec} + \lambda \mathcal{L}_{cl}, \quad (1)$$

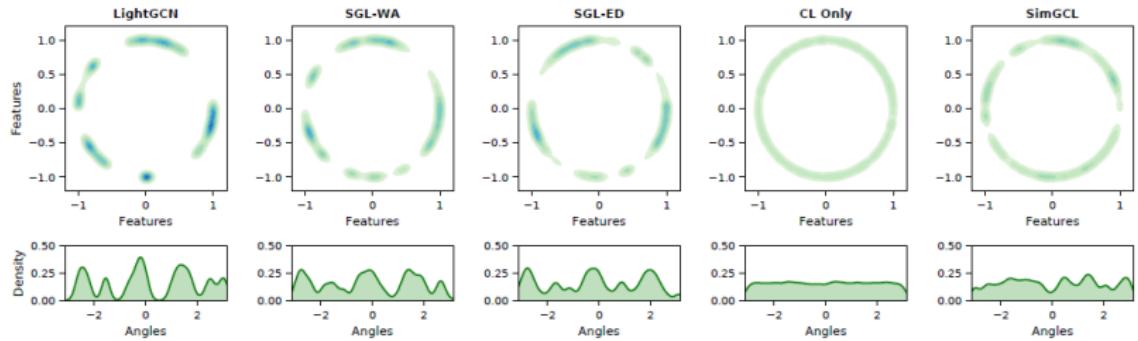
which consists of two losses: recommendation loss \mathcal{L}_{rec} and CL loss \mathcal{L}_{cl} . The InfoNCE in SGL is formulated as:

$$\mathcal{L}_{cl} = \sum_{i \in \mathcal{B}} -\log \frac{\exp(\mathbf{z}_i'^\top \mathbf{z}_i'' / \tau)}{\sum_{j \in \mathcal{B}} \exp(\mathbf{z}_i'^\top \mathbf{z}_j'' / \tau)}, \quad (2)$$

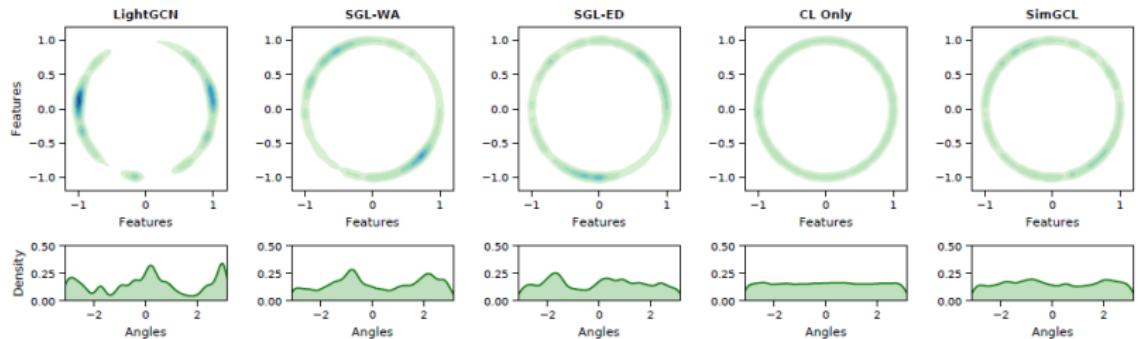
- L_2 normalized embeddings
- differences between contrastive views
- the key factor in SGL was the loss function (not the augmentations with edge/node dropout)

InfoNCE: <https://arxiv.org/abs/1807.03748>

Embedding Uniformity



(a) Distribution of item representations learned from the dataset of Yelp2018.



(b) Distribution of item representations learned from the dataset of Amazon-Book.

[uncovered] Augmentations for Interactions

[uncovered] DiffuASR

Diffusion

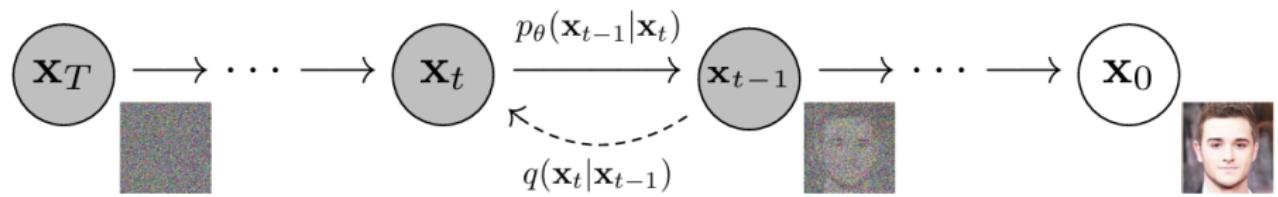


Figure 2: The directed graphical model considered in this work.

DiffuASR

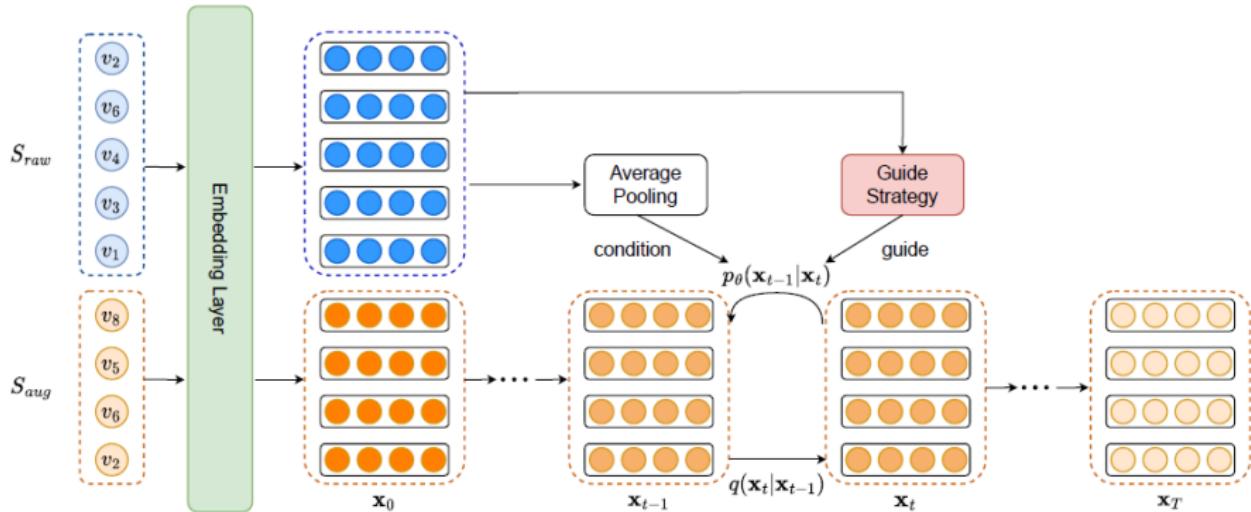


Figure 3: The overview of the proposed Diffusion Augmentation for Sequential Recommendation (DiffuASR).

<https://dl.acm.org/doi/10.1145/3583780.3615134>

SU-Net

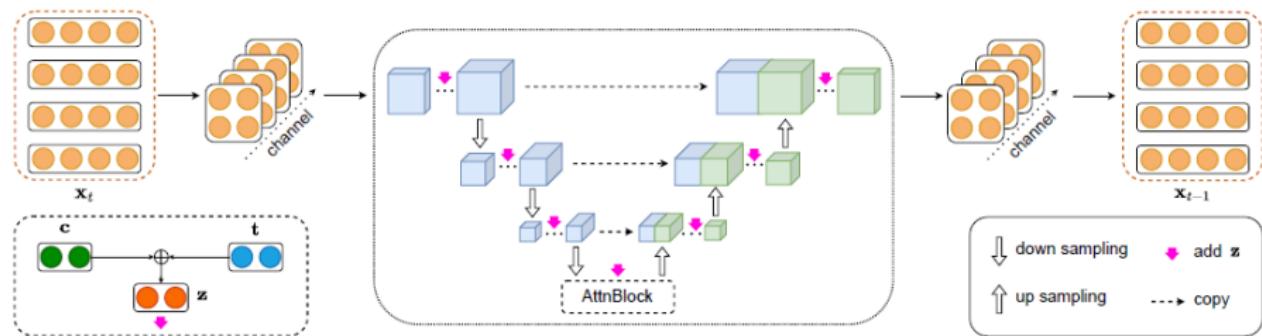


Figure 4: The architecture of the proposed SU-Net.

[uncovered] Matrix Form – revisited

[uncovered] Matrix Form – revisited

[uncovered] DiffRec

DiffRec

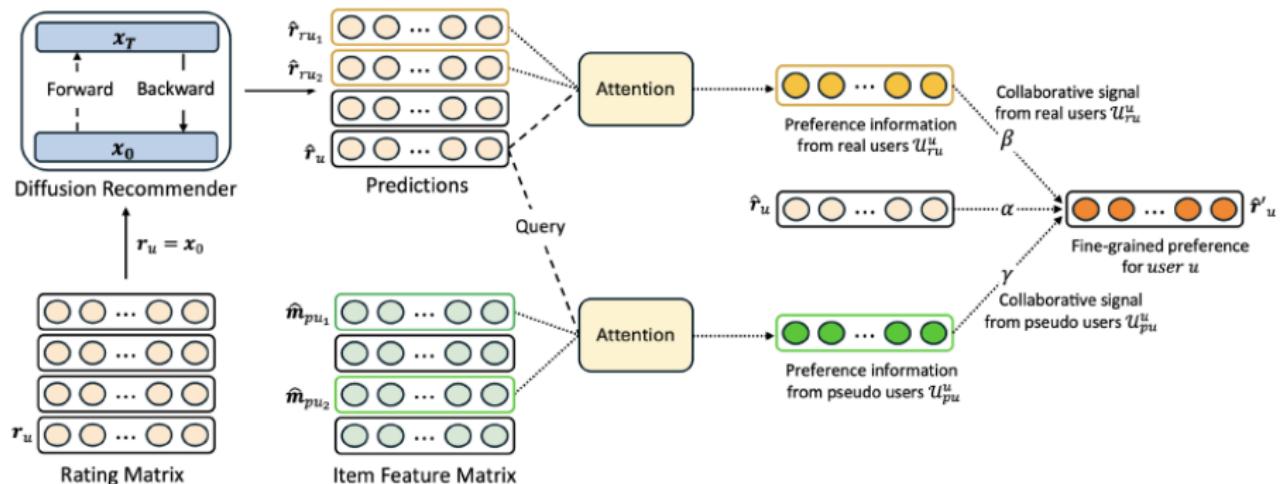


Figure 1: The overview of the proposed Collaborative Diffusion Model for Recommender System (CDiff4Rec).

<https://arxiv.org/abs/2501.18997>