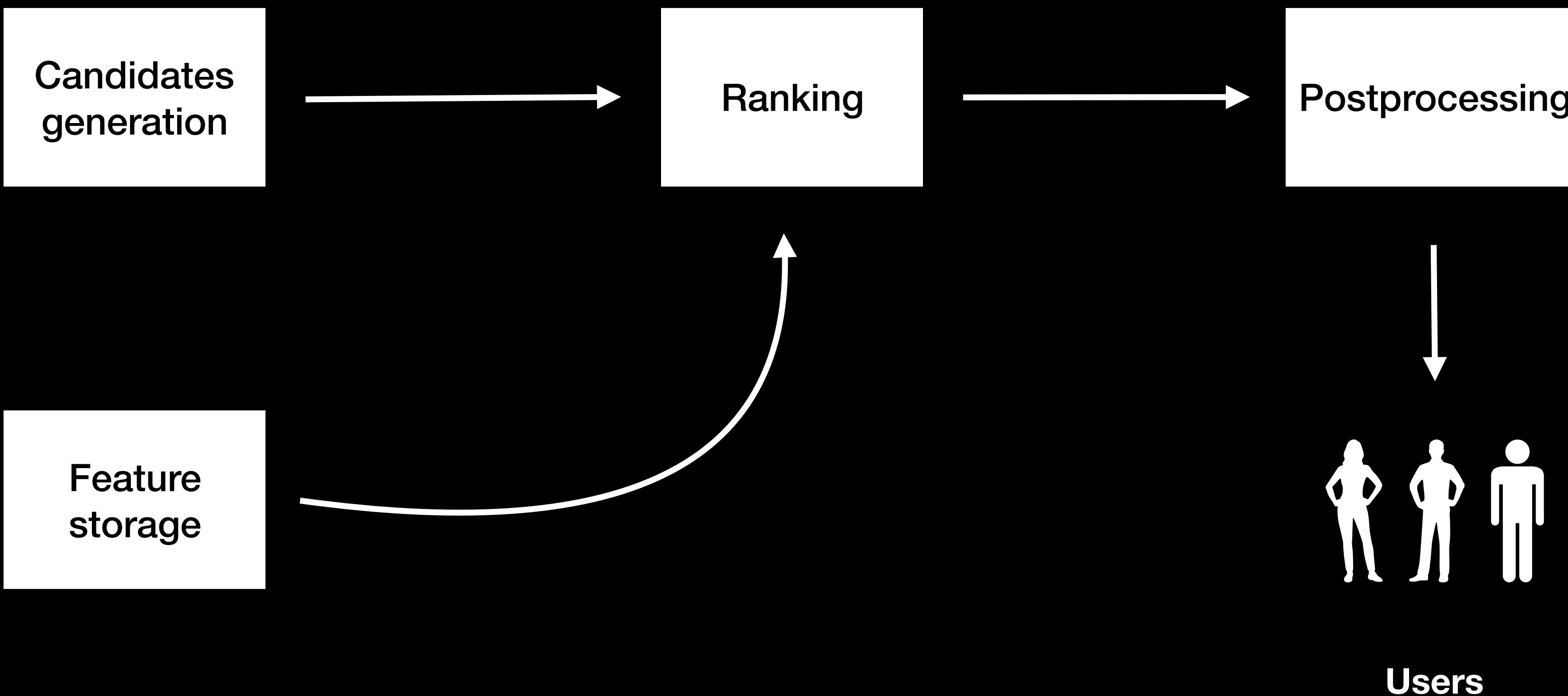


# Recommendation Systems

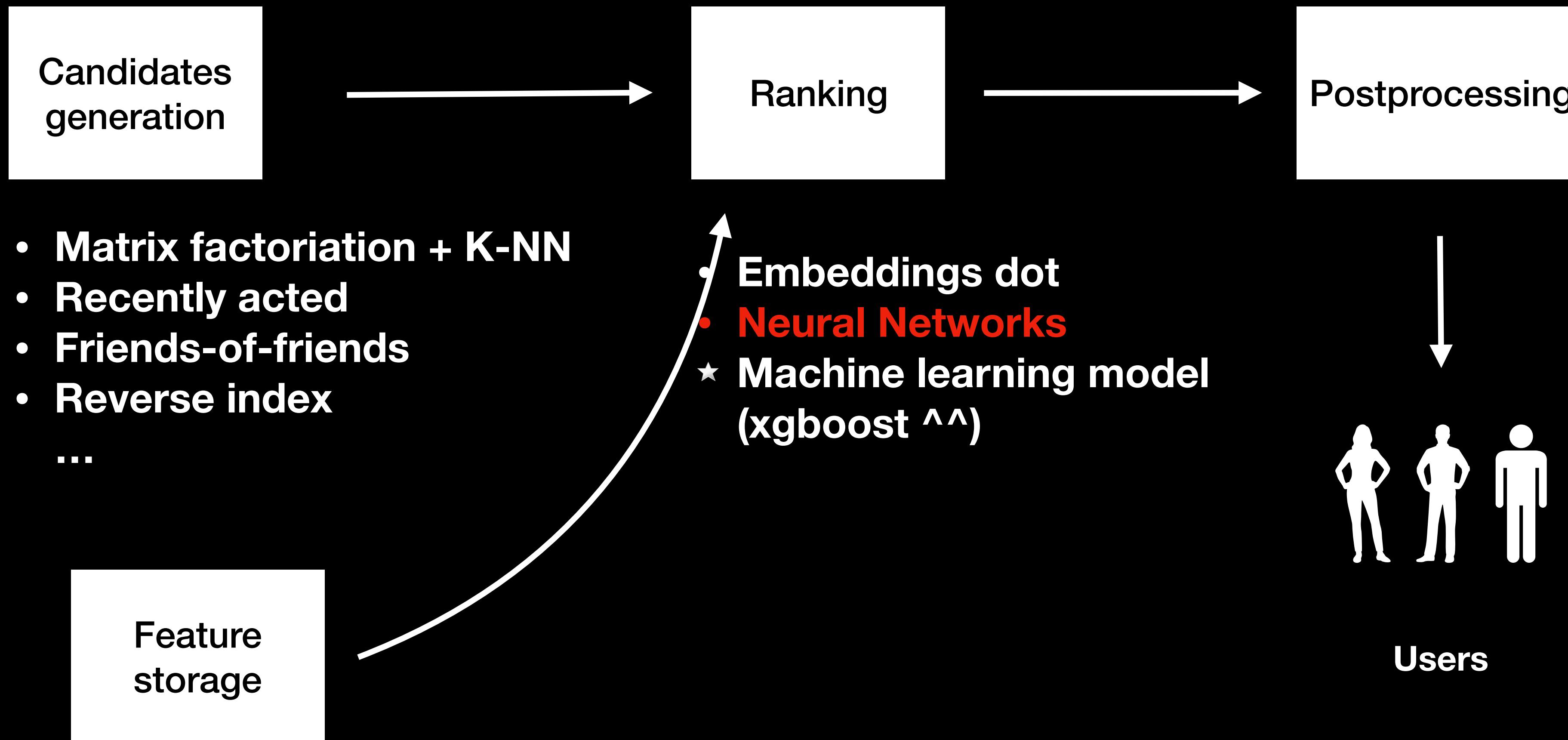
**Deep Learning, part 2**

Eugeny Malyutin / Sergey Dudorov

# Recommendations pipeline:



# Recommendations pipeline:



# Neural networks as ranking:

- Task: airbnb search.

Stop, what? We are about recsys. Are you **mad?** Nope. Have you seen airbnb search? It's totally about ranking and search is only «name» , no normal textual questions, all query is built from simple geo and price ranges so yes, it's **more recommender system then search.** They use ranking objectives, measures NDCG and got strong correlations between NDCG- and bookings- gains (offline-online correlation) that's why it's more recommendation then search

- Dataset: booked vs viewed.  
Offline quality: NDCG.  
Baseline: GBDT.  
Online metric: bookings.



# Ch1. Do not be a heroes:

Custom super-designed architecture.

The same GBDT features.

L2-loss with 1 as positives, 0 as negatives.

Awful results.

**WHY?**

# Ch1. Do not be a heroes:

Simple NN with 32-units hidden layer.

The same GBDT features.

L2-loss with 1 as positives, 0 as negatives.

Awful results.

**WHY?**

# Ch1. Do not be a heroes:

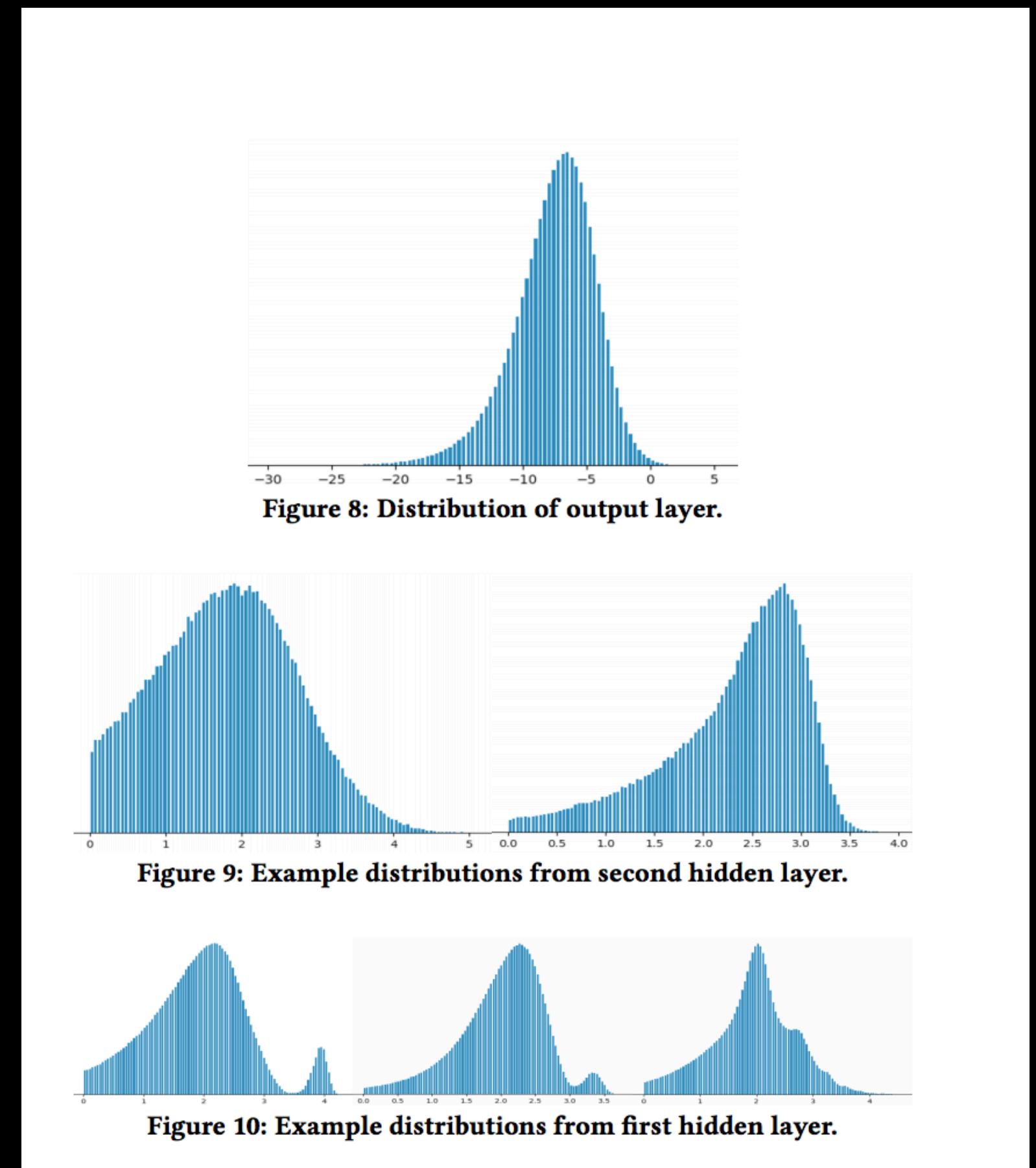
Simple NN with 32-units hidden layer. The same GBDT features. L2-loss with 1 as positives, 0 as negatives. Awful results.

- Feature normalisation:
  - Normal distribution:  $(f\_val - \mu)/\sigma$
  - Power law distr:  $\log(\frac{1 + f\_val}{1 + median})$

# Ch1. Do not be a heroes:

Simple NN with 32-units hidden layer. The same GBDT features. L2-loss with 1 as positives, 0 as negatives. Awful results.

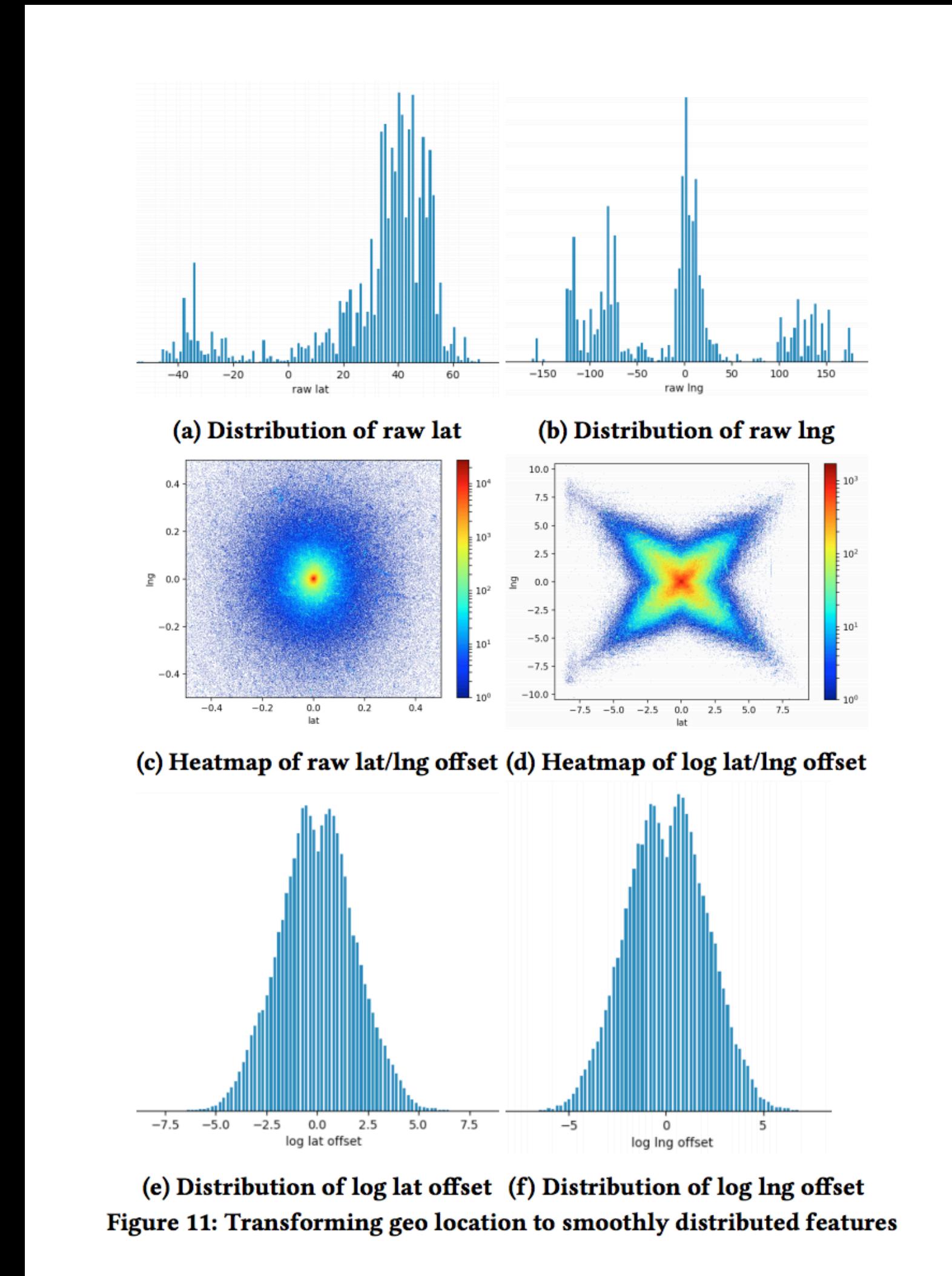
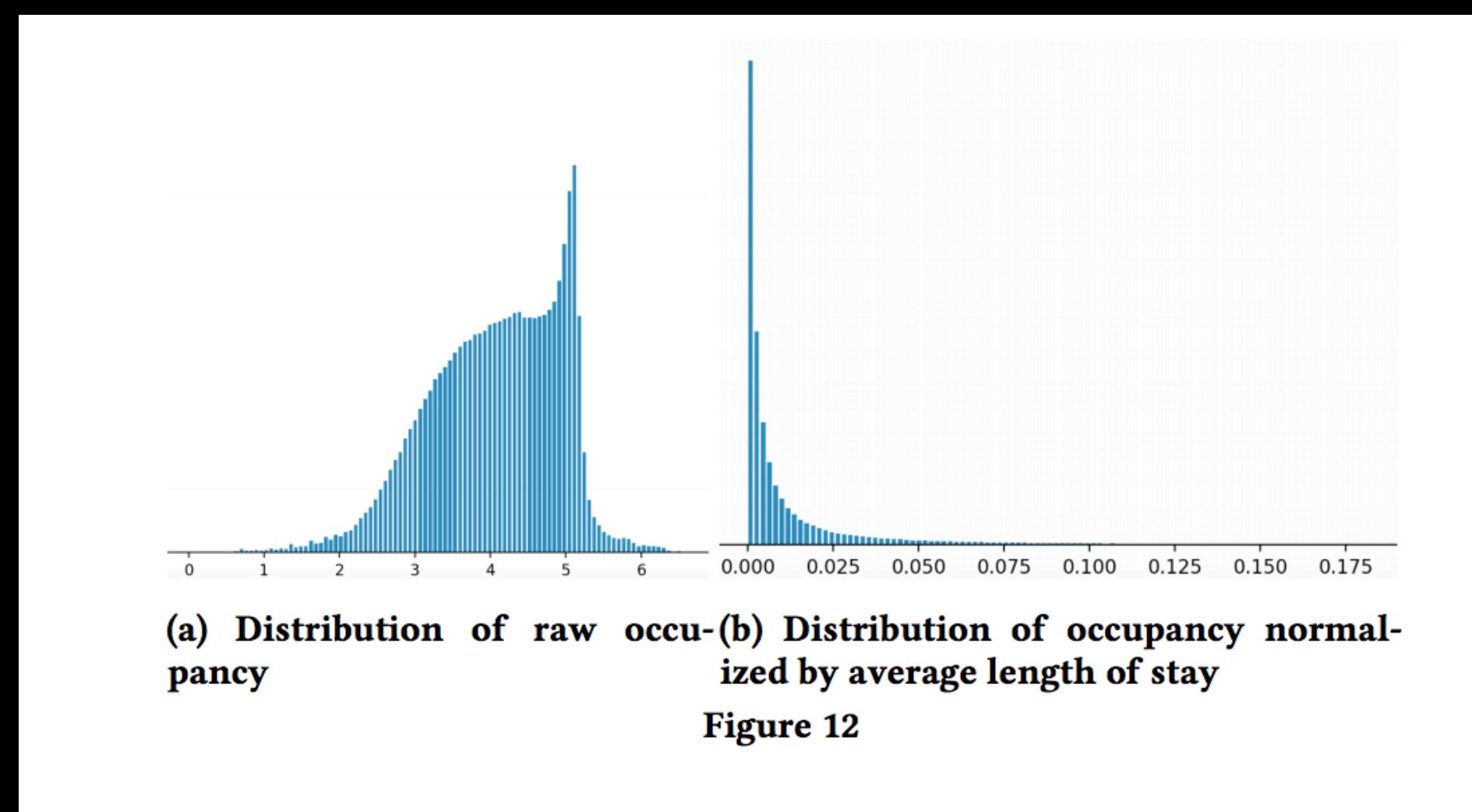
- Smooth distributions:
  - Spotting bugs
  - Facilitating generalization



# Ch1. Do not be a heroes:

Simple NN with 32-units hidden layer. The same GBDT features. L2-loss with 1 as positives, 0 as negatives. Awful results.

- Special features mapping. Lat/lng to offset wrt users screen centered.
- Completeness of features: occupancy



# Ch1. Do not be a heroes:

Simple NN with 32-units hidden layer.

**Features are normalized and validated.**

L2-loss with 1 as positives, 0 as negatives.

Normal results.

# Ch2. Main character runs up a stairway:

LambdaRank in NN:

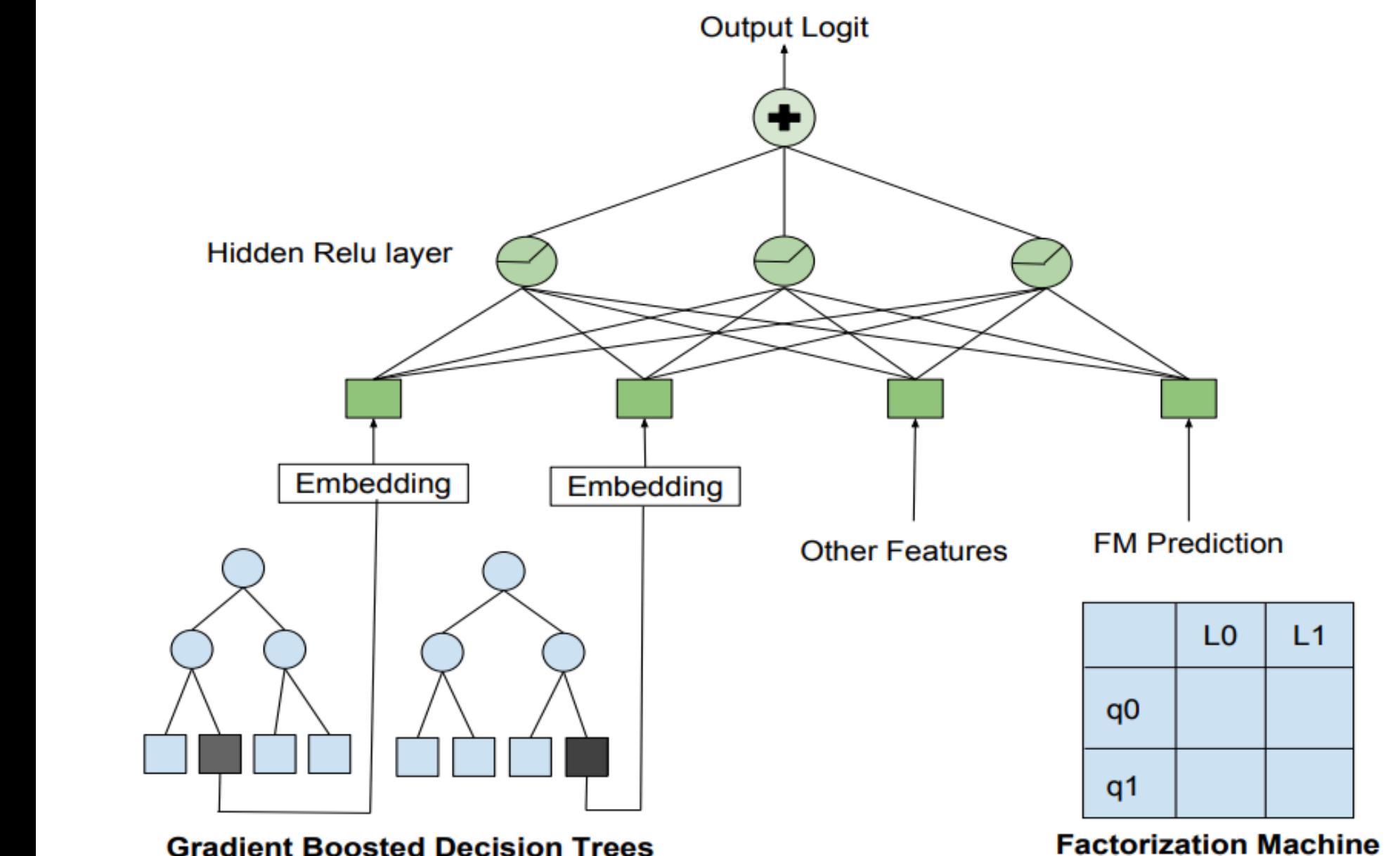
- Pairwise dataset. {positive\_list, negative\_list} sampled from the same session logs. CE loss between listed and not-listed
- Scaled at NDCG change of swapping pairs.
- ^^ Prioritizes 1 to 0 ranking over 125 to 124.

# Ch2. Main character runs up a stairway:

GBDT-NN-shtein:

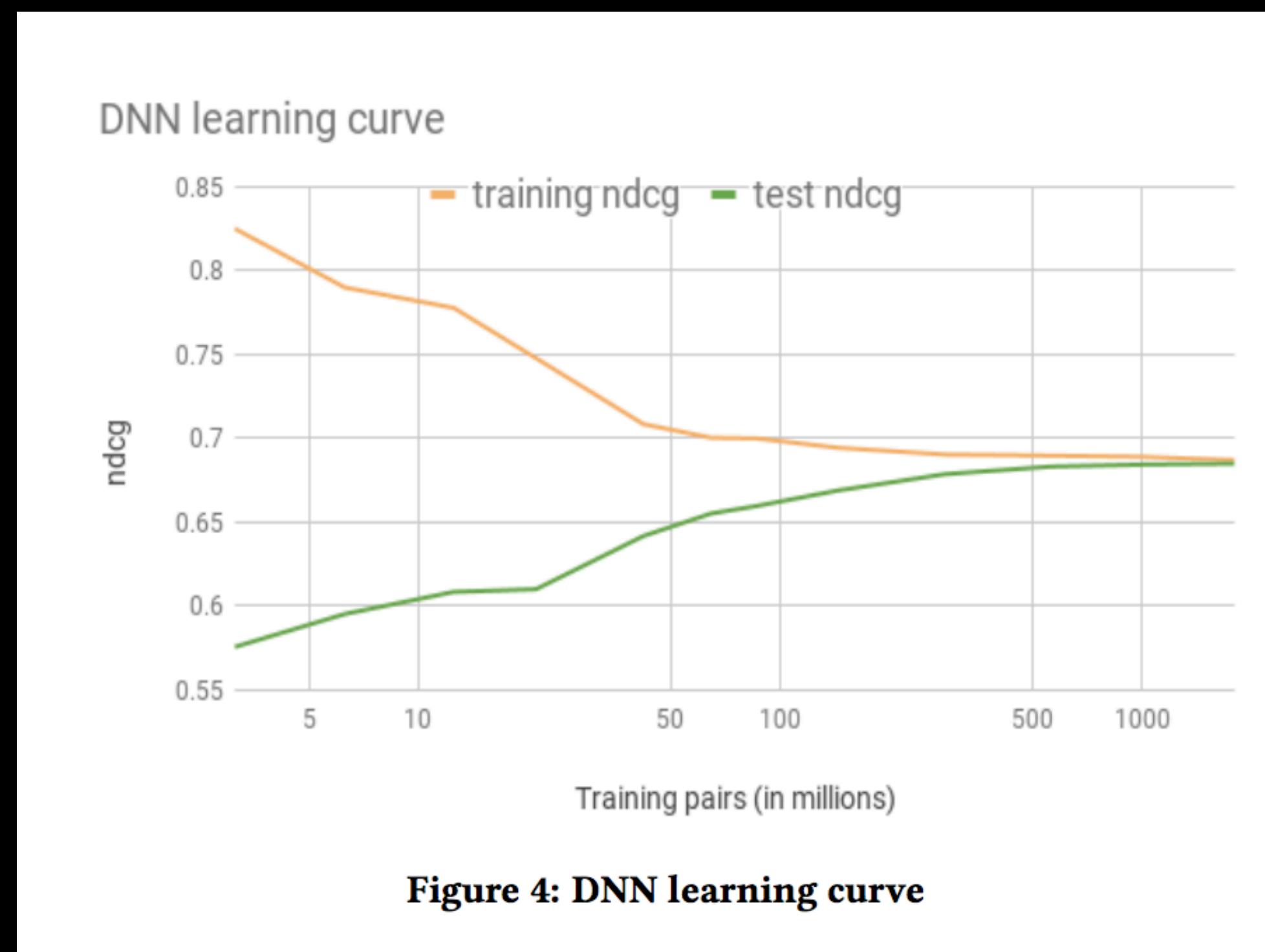
- GBDT up-ranks different samples.
- FM-machines with their embeddings.
- NN with pretty smooth generalisation and/or interactions

Frankenstein to combine them all ->



# Ch2. Main character runs up a stairway:

Frankenstein to combine them all -> was overperformed by this network and 10x data.



LambdaNN

97 RELU

127 RELU

195 features

# Ch2. Main character runs up a stairway:

- «General» embeddings (but location preferences)

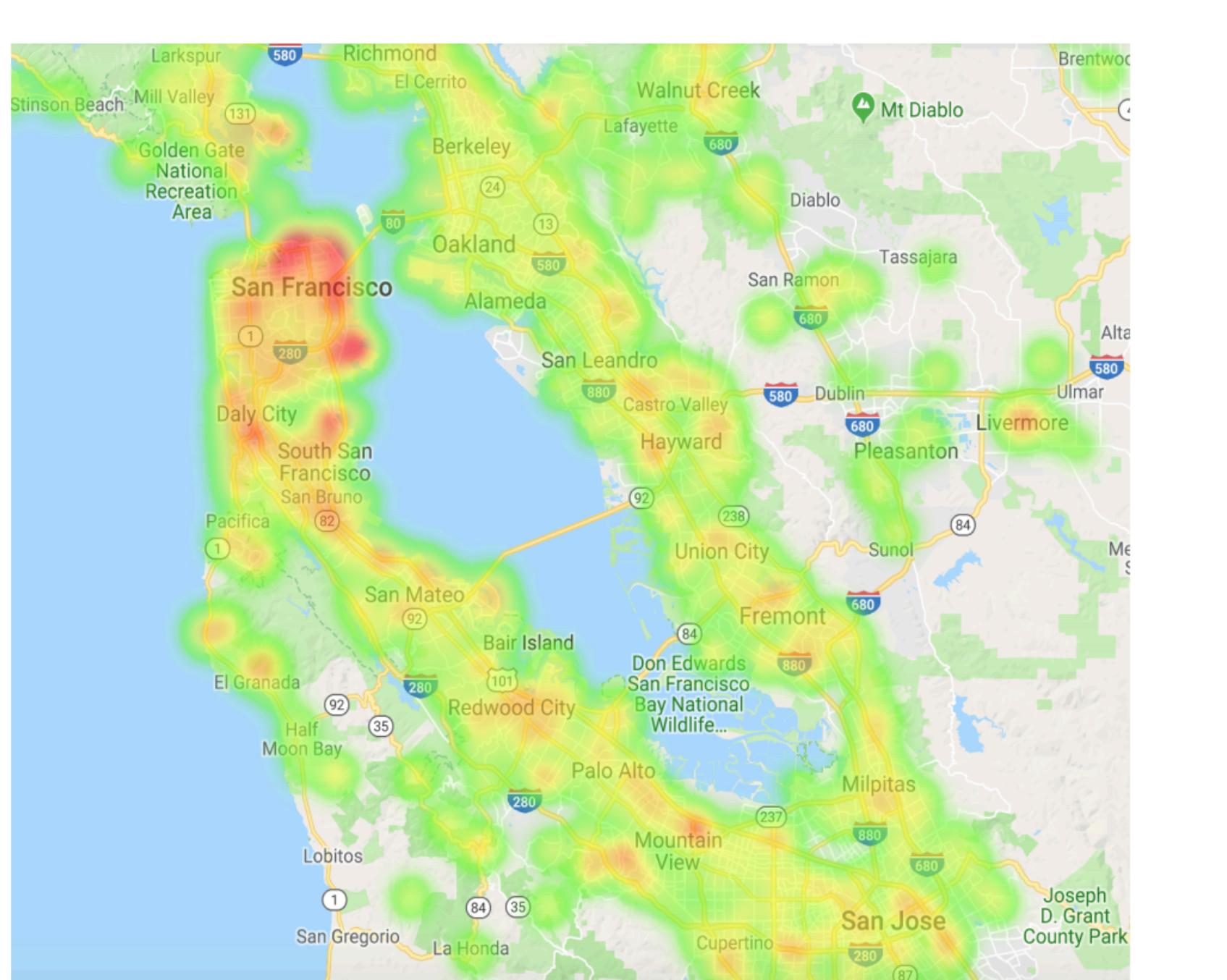


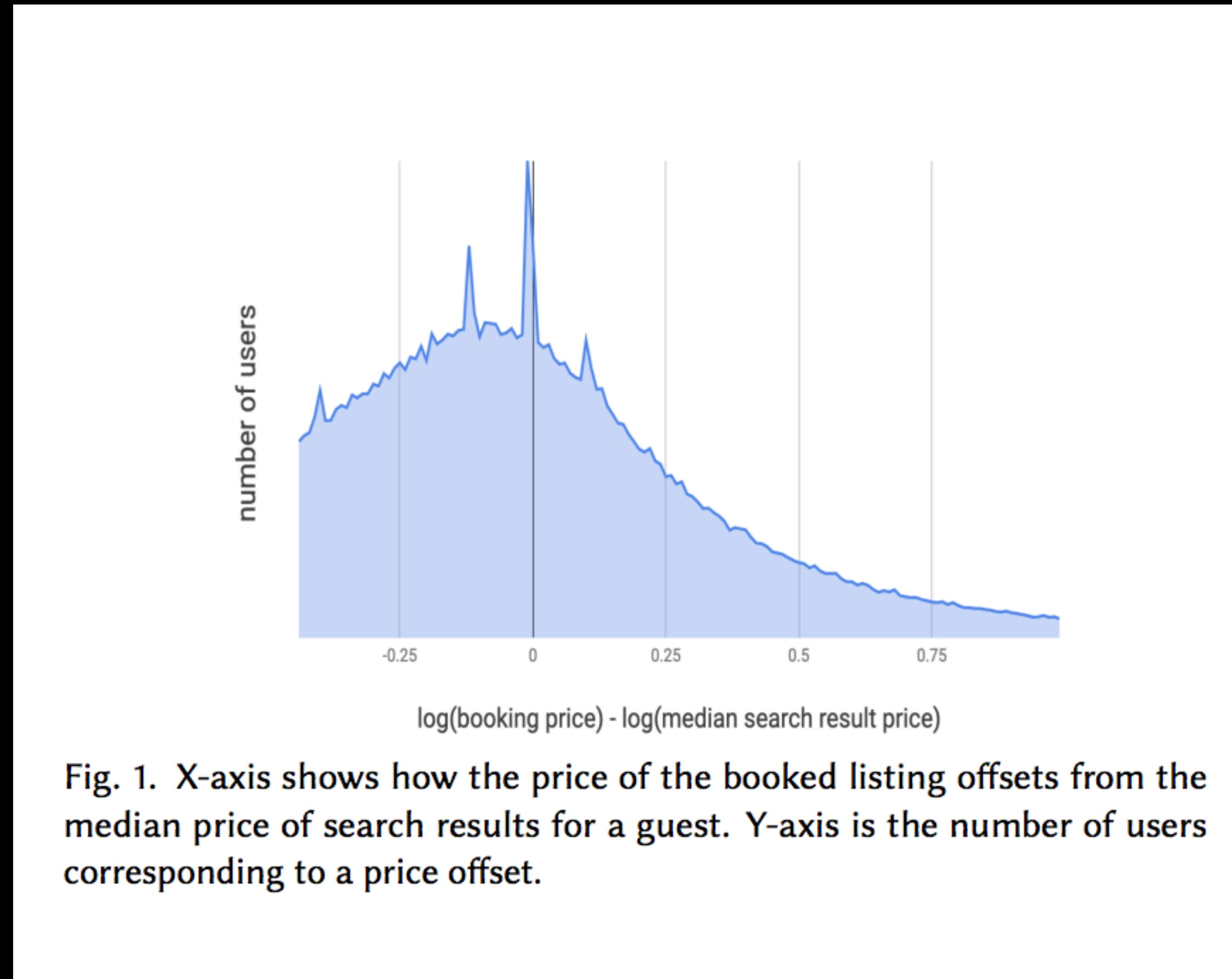
Figure 13: Location preference learnt for the query "San Francisco"

- {"San Francisco", 539058204} → 71829521

# Ch3. Hero recall failed brothers:

- MOAR layers
- «Modern attention architectures»
- Direct entities embeddings
- ...

# Vol.2 Hero strikes back



# Ch1: Naive

$$DNN_{\theta}(u, q, l_{noPrice}) + (-\tanh(w * P + b))$$

$$P = \log\left(\frac{1 + price}{1 + price_{median}}\right)$$

- -5.7% avg price in booking, -1.5% of bookings

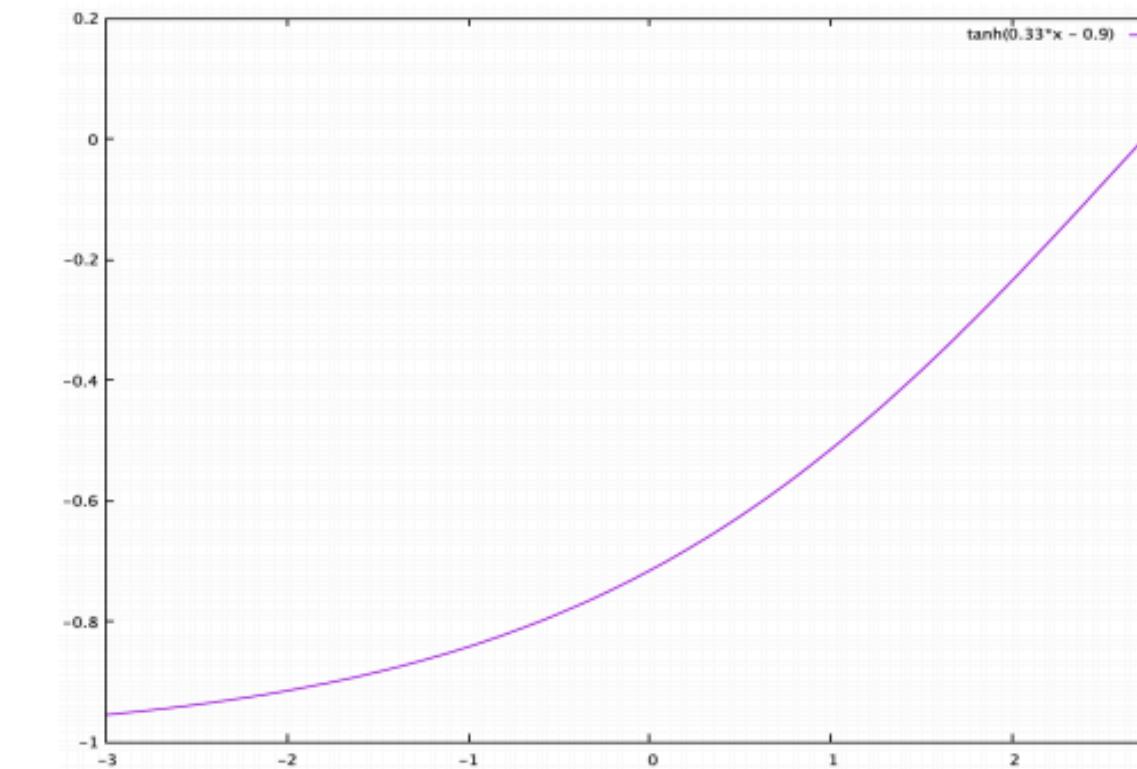


Fig. 2. X-axis is normalized price feature. Y-axis is the value of the  $\tanh$  term in Equation 1.

# Ch1: Less naive

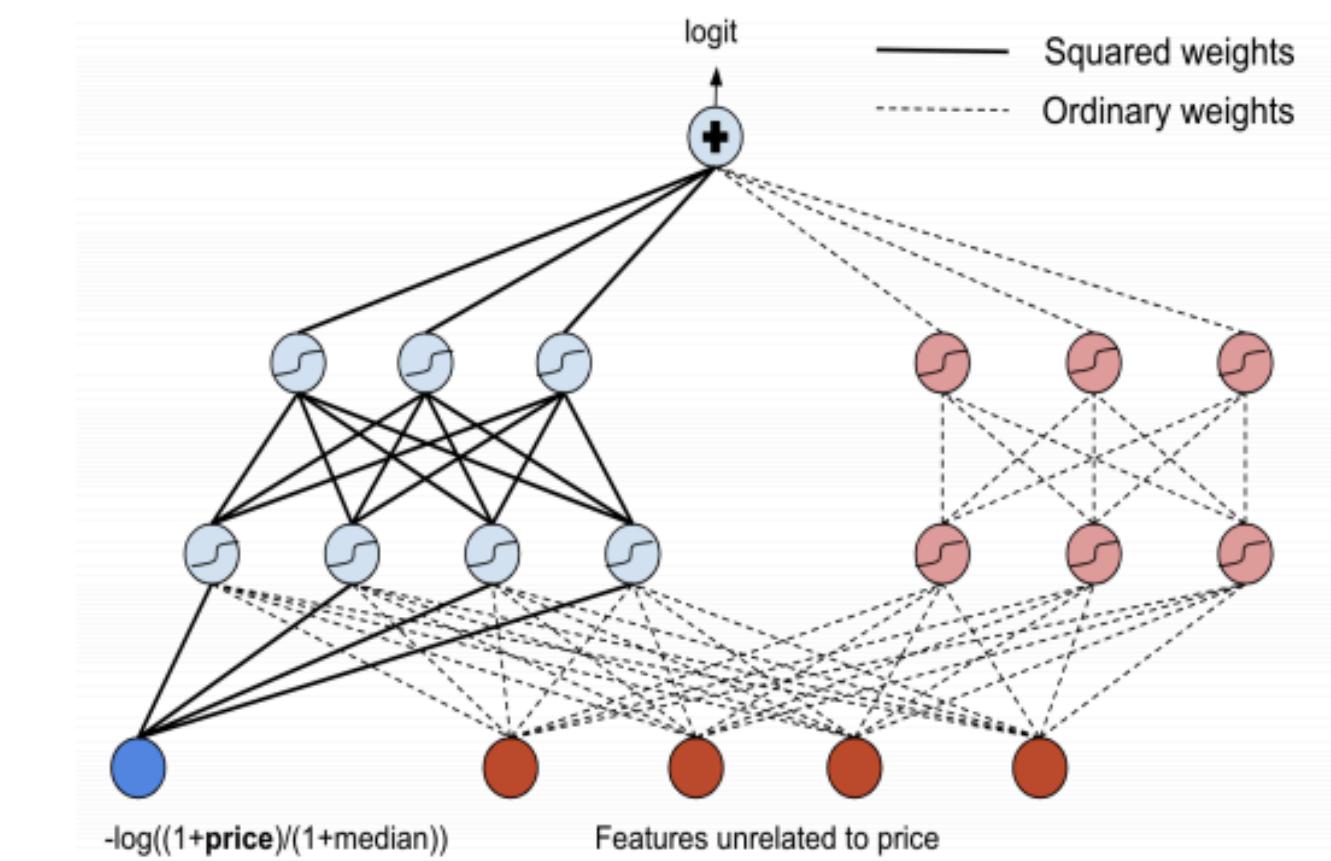


Fig. 3. DNN architecture partially monotonic w.r.t price. Bold solid lines indicate weights that are squared, dashed lines indicate ordinary weights.

# Ch1: Less naive

- **-1.4% of bookings**
- Tried even additional loss component - **nothing.**

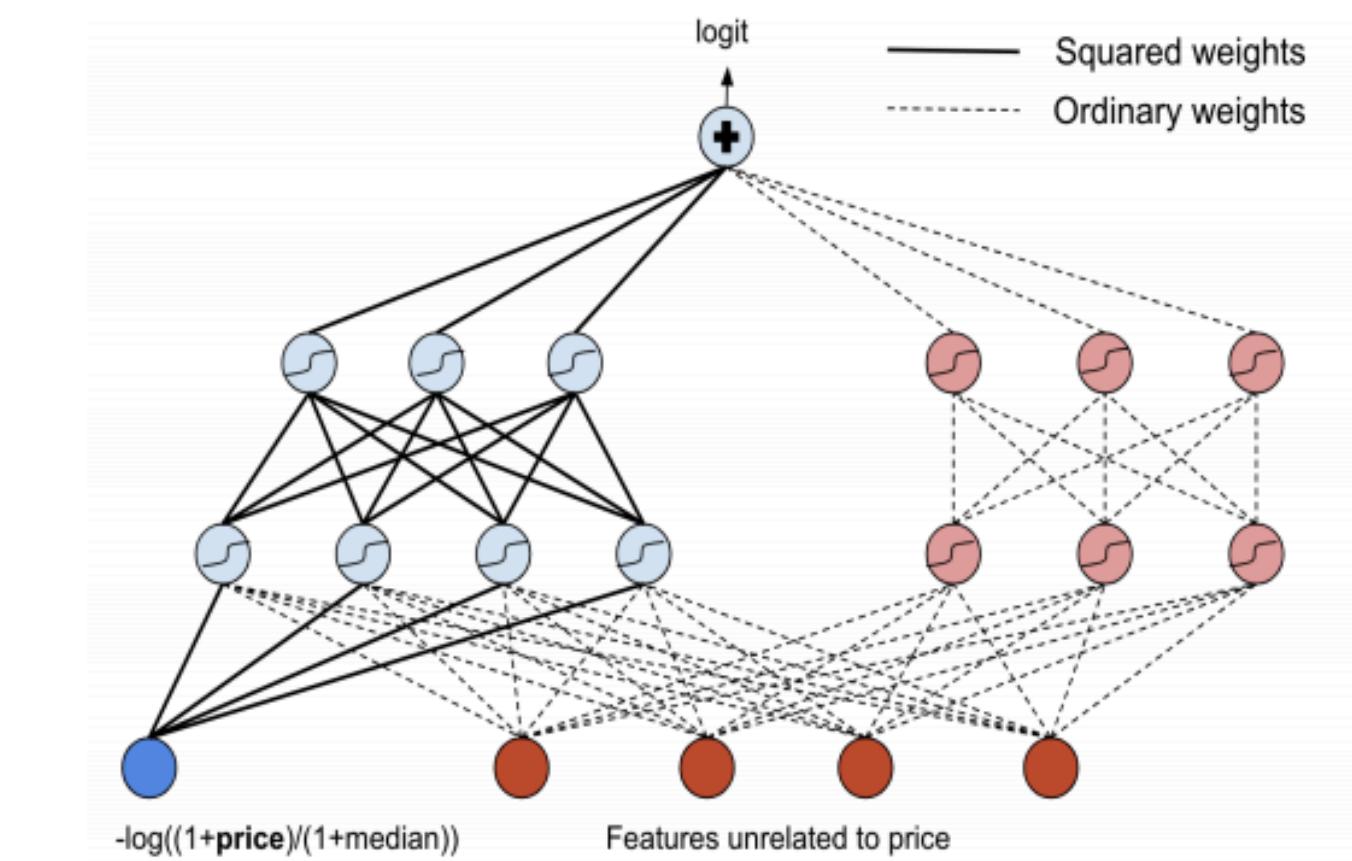


Fig. 3. DNN architecture partially monotonic w.r.t price. Bold solid lines indicate weights that are squared, dashed lines indicate ordinary weights.

# Ch2: I - interpretability

- Individual conditional expectation (ICE) plots

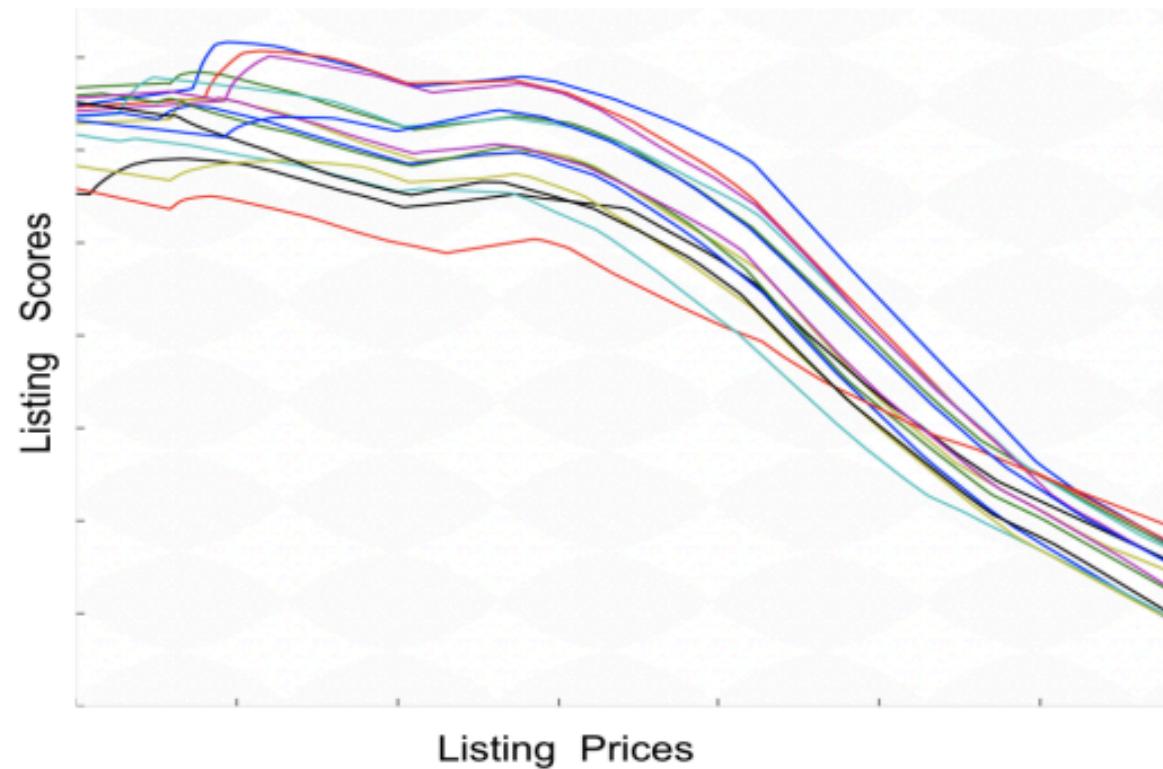
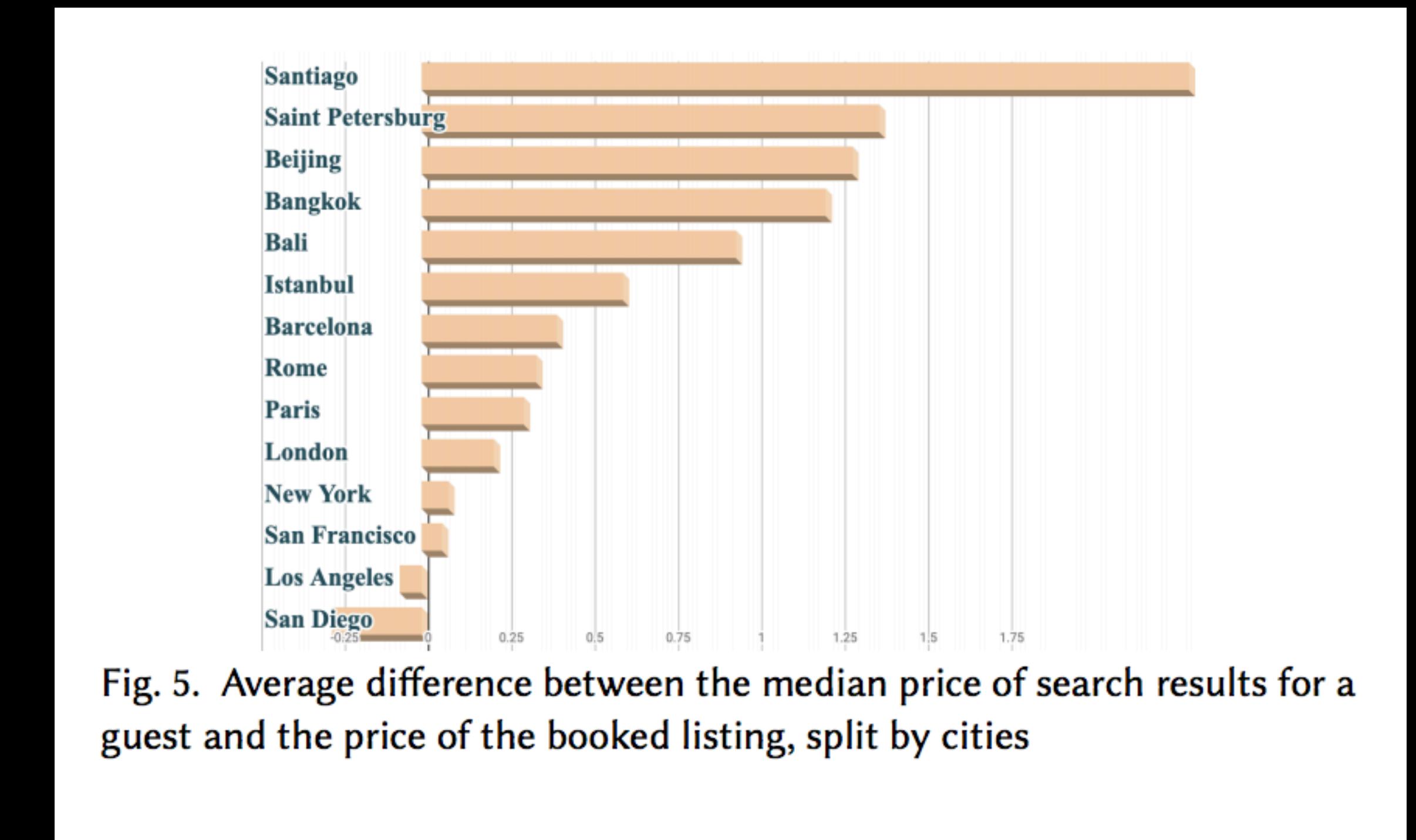


Fig. 4. ICE plot of listing scores (y-axis) vs listing prices (x-axis). Each curve represents a listing in the search result.

# Ch2: I - interpretability

- Individual conditional expectation (ICE) plots



# Ch2: I - interpretability

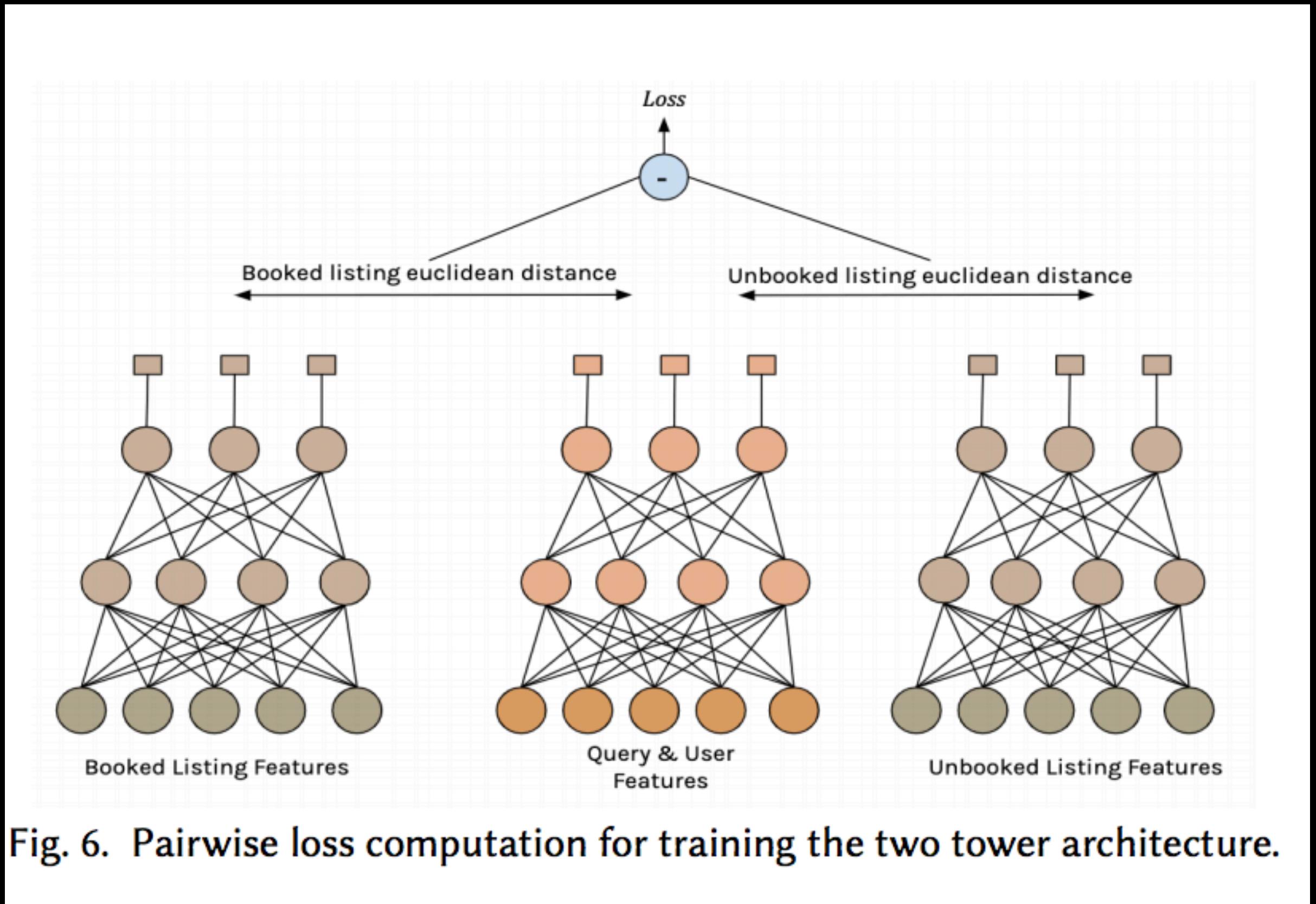


Fig. 6. Pairwise loss computation for training the two tower architecture.

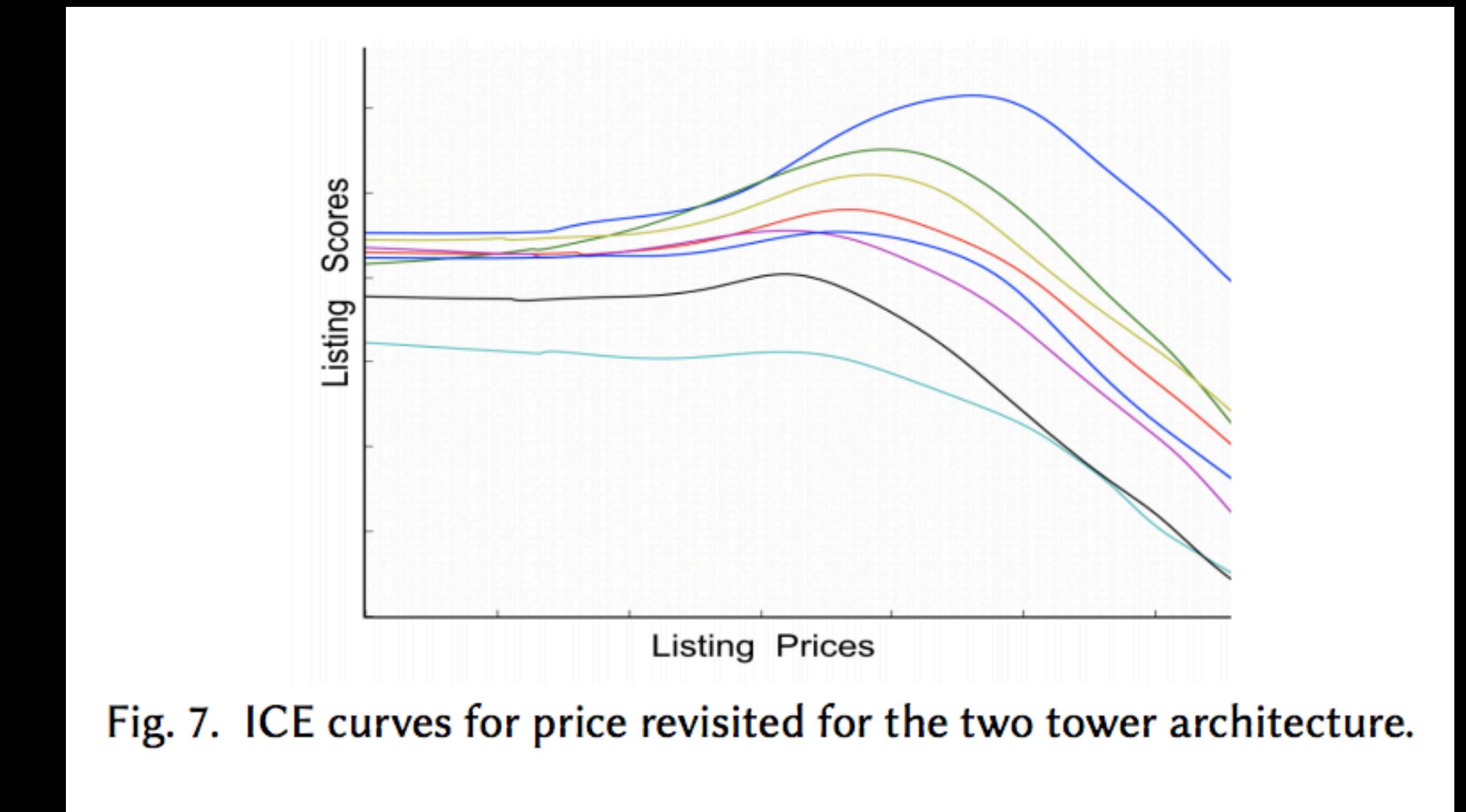
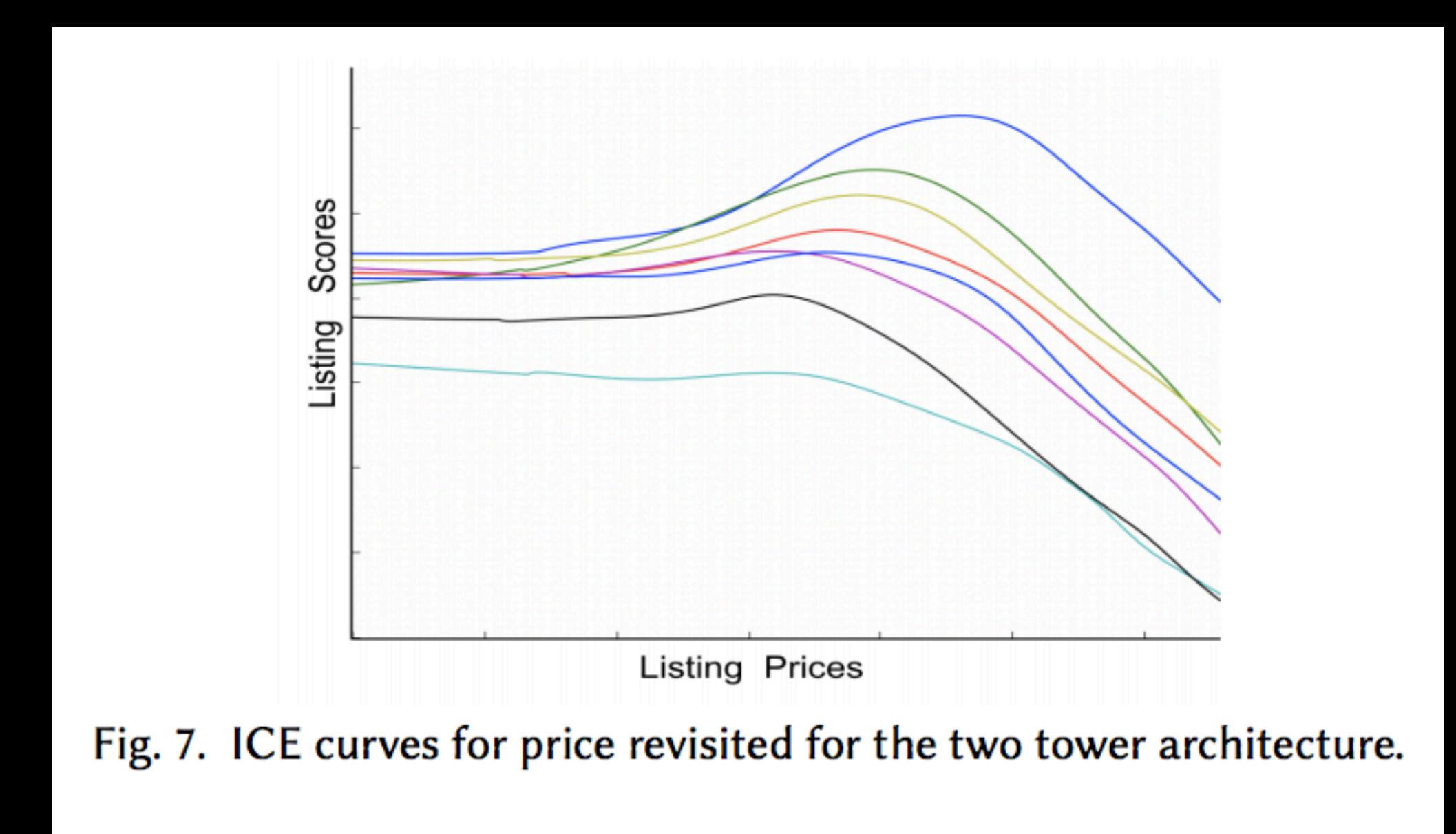


Fig. 7. ICE curves for price revisited for the two tower architecture.

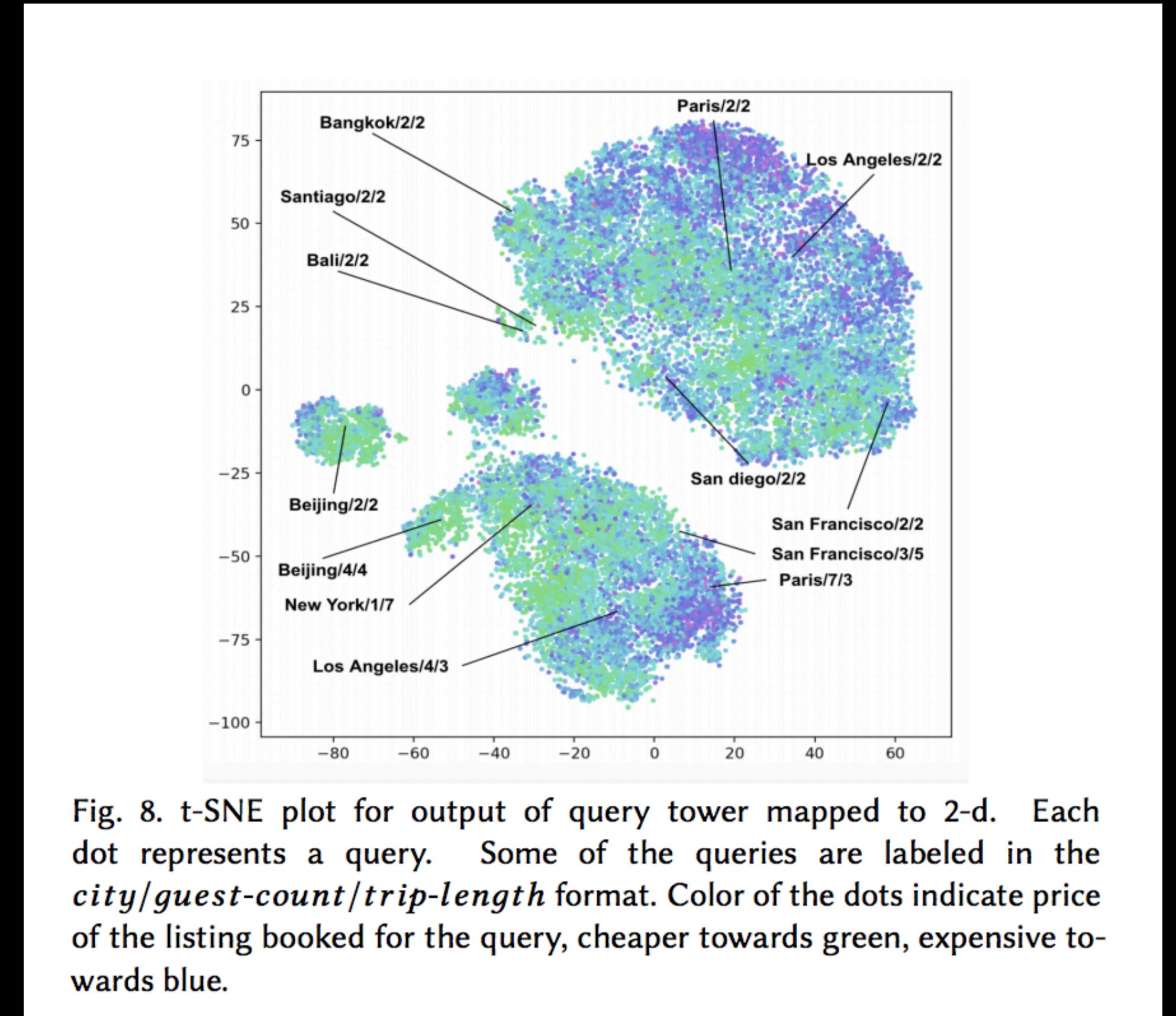
- **-2.3% of bookings, +0.75% revenue**
- -33% peak latency(!)

# Ch2: I - interpretability



# Ch3: Lessons learned

- Users lead - model follows.
- Modern architectures highly tied with applications
- Architecture engineering, not complications.
- Read moar:
  - Cold start
  - Positional bias
  - Improving Deep Learning For Airbnb Search (c)



# Modern architectures are like:

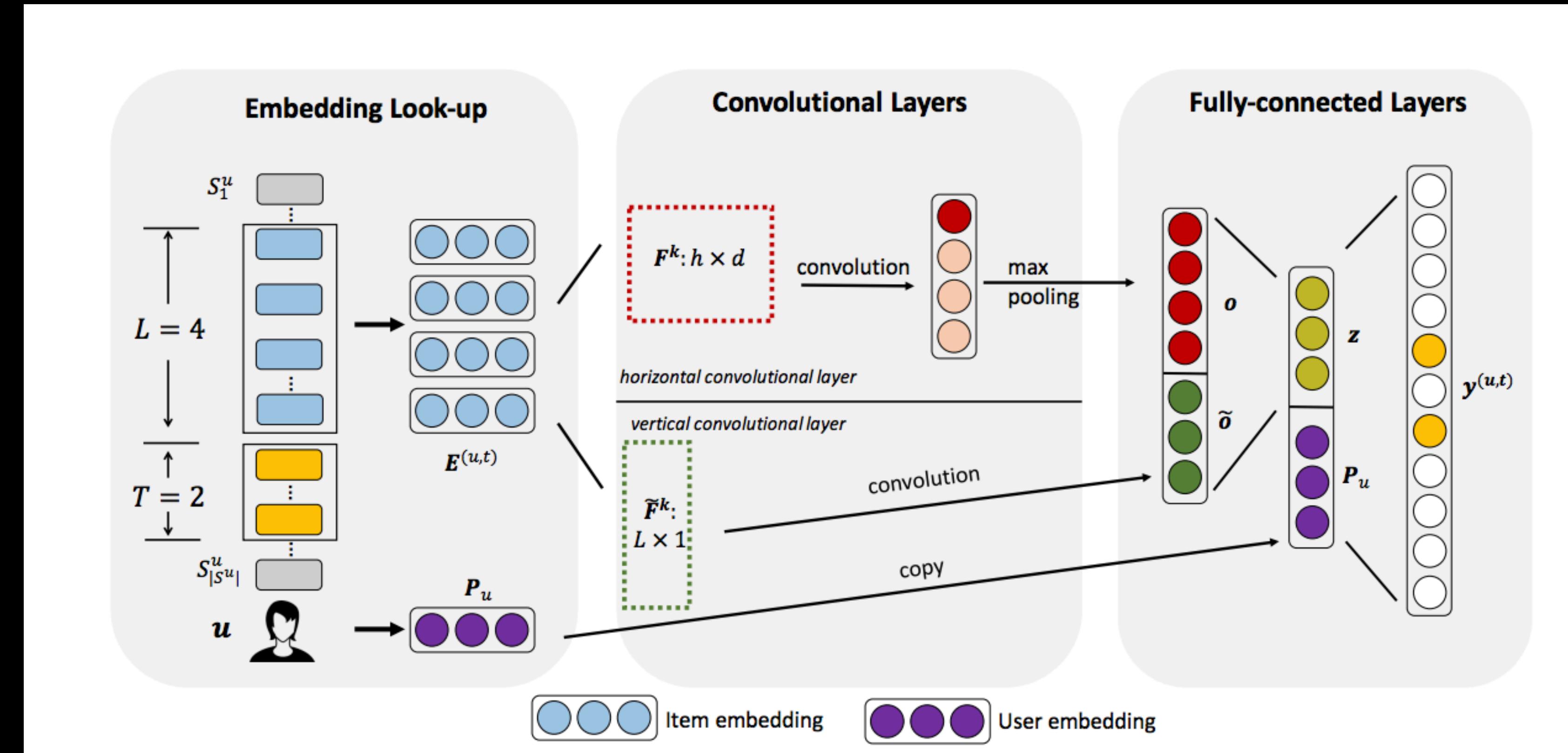


RecSys

The rest of DL

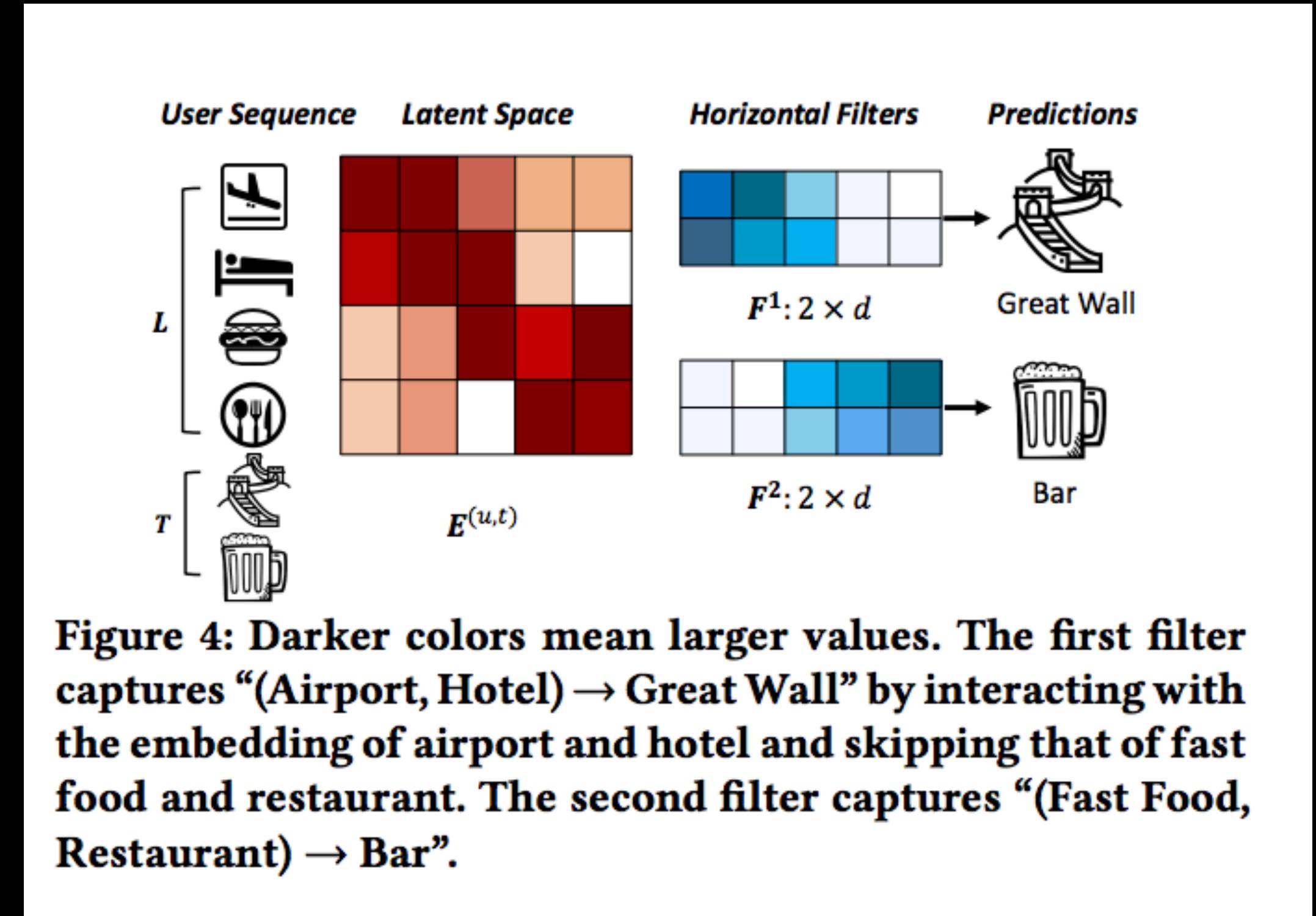
# Caser:

- $L$  – history size
- $T$  – num of steps to predict
- $S_t^u$  – user interacted item at time
- $P_u$  – user embedding



# Caser:

- Horizontal convolution - to capture «semantic» pattern. Capture union-level patterns with multiple union sizes (c)
- These vertical Iters are capturing point-level sequential patterns through weighted sums over previous items' latent representations. (c)
- Binary cross-entropy loss for next ( $T(!)$ ) items to allow «skip-behaviour».



**Figure 4: Darker colors mean larger values. The first filter captures “(Airport, Hotel) → Great Wall” by interacting with the embedding of airport and hotel and skipping that of fast food and restaurant. The second filter captures “(Fast Food, Restaurant) → Bar”.**

What's went wrong?

# Caser:

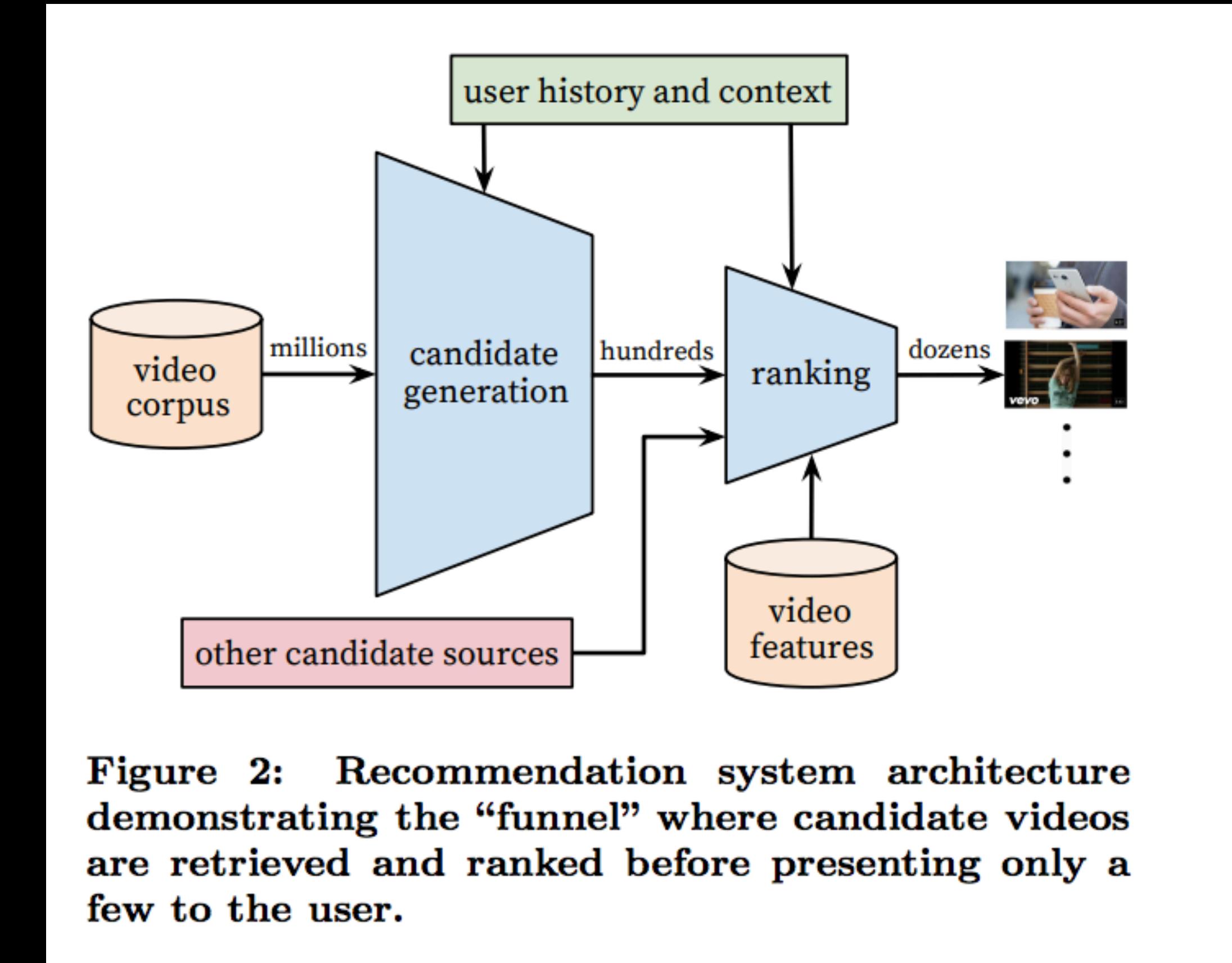
- Horizontal convolution - to capture «semantic» pattern. Capture union-level patterns with multiple union sizes (c)
- These vertical Iters are capturing point-level sequential patterns through weighted sums over previous items' latent representations. (c)
- Binary cross-entropy loss for next ( $T(!)$ ) items to allow «skip-behaviour».
- Lack of capacity limited to size of a filters.
- Factually  $E(u,t)$  matrix has no geo-spatial local feature interactions. Learning such embeddings factually trickier then classic one.
- Binary CE...:
  - How to add new items?
  - How to add new users?
  - How can we re-evaluate network with new user activity?

What's went wrong?

# YT(youtube) recommender:

- Loss:  $P(w_t = i \mid U, C) = \frac{e^{v_i, u}}{\sum_{j \in J} e^{v_j, u}}$

Problems?



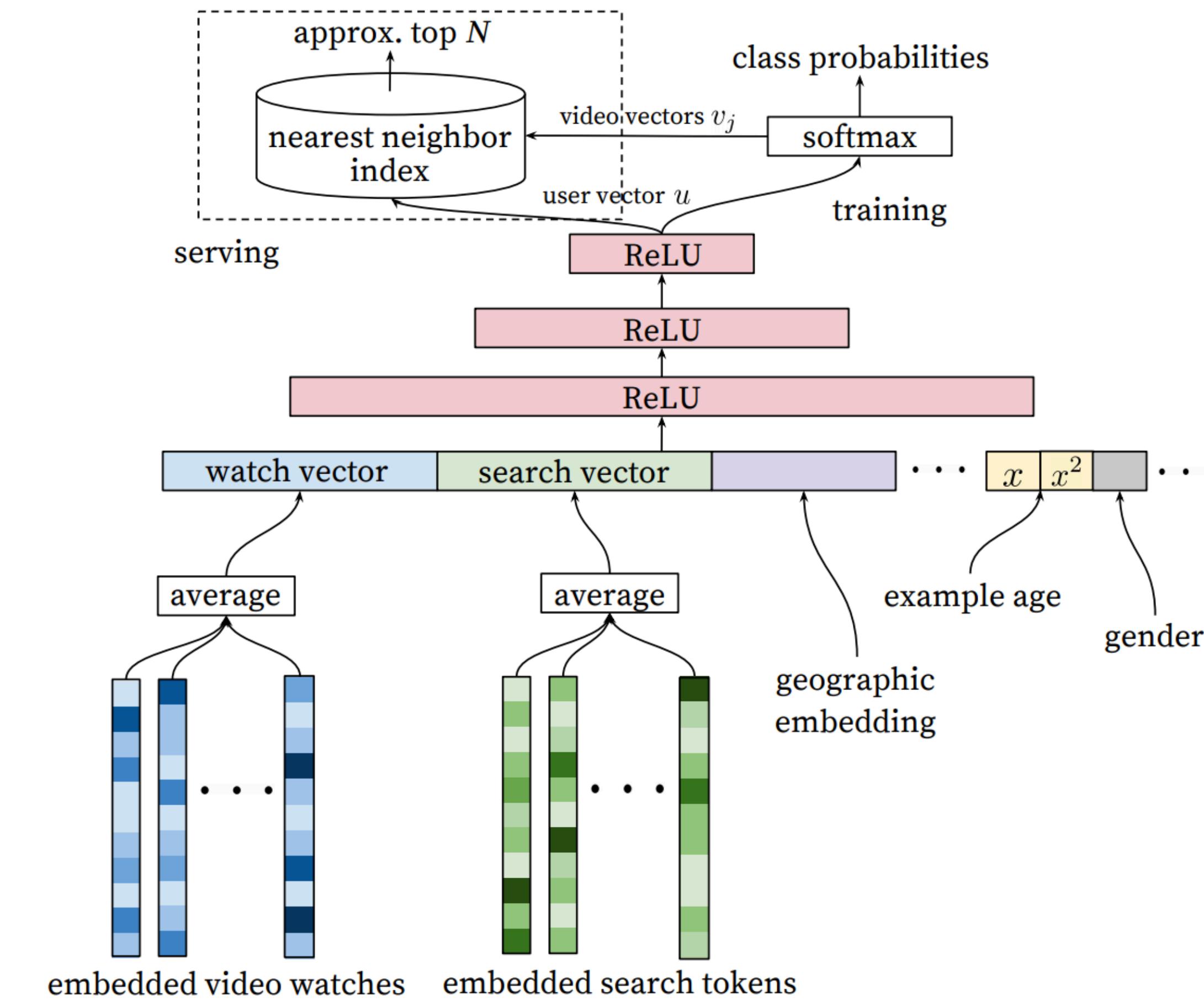
Source:

# YT(youtube) recommender:

- Loss:

$$P(w_t = i | U, C) = \frac{e^{v_i, u}}{\sum_{j \in J} e^{v_j, u}}$$

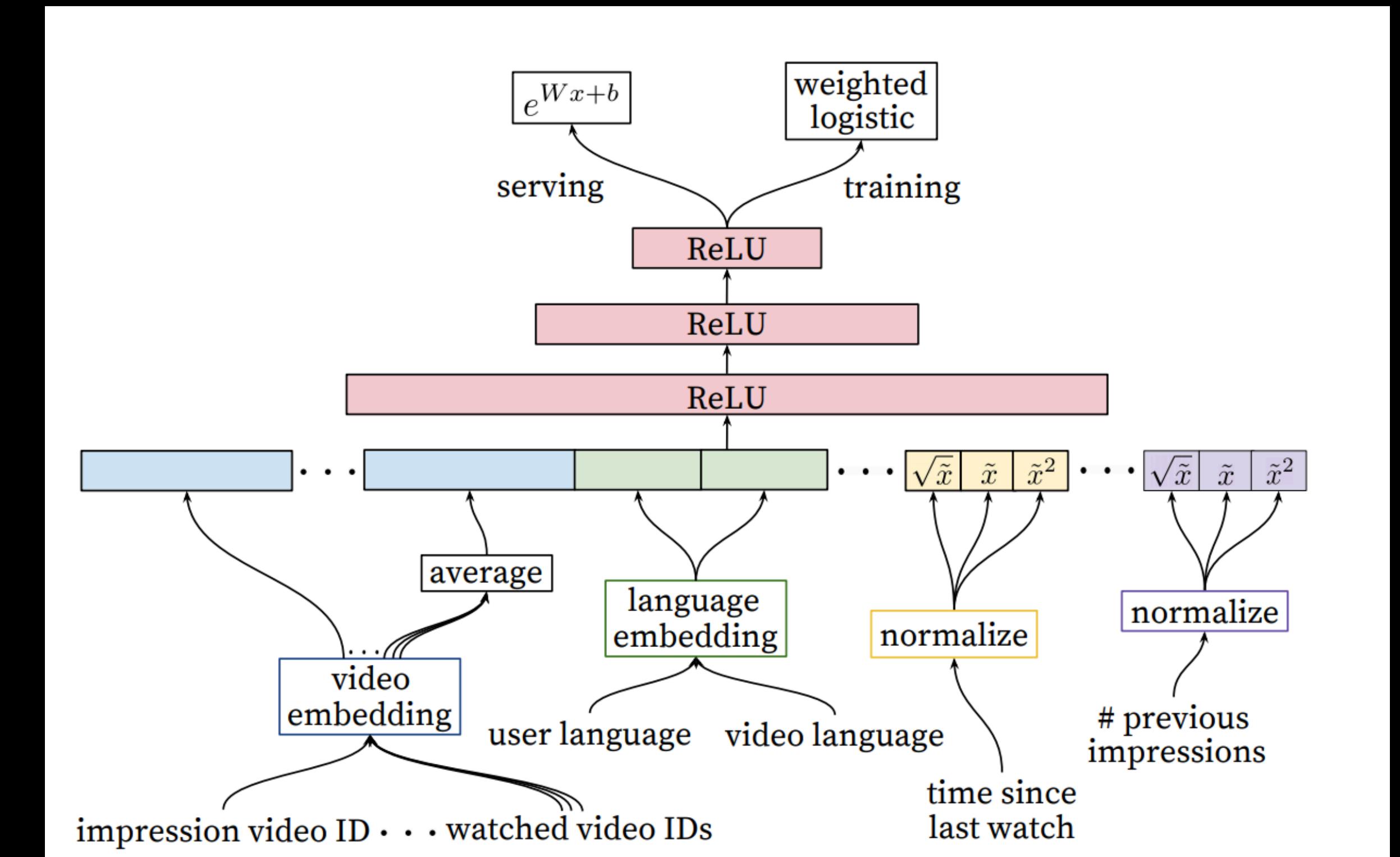
Yep, neg sampling



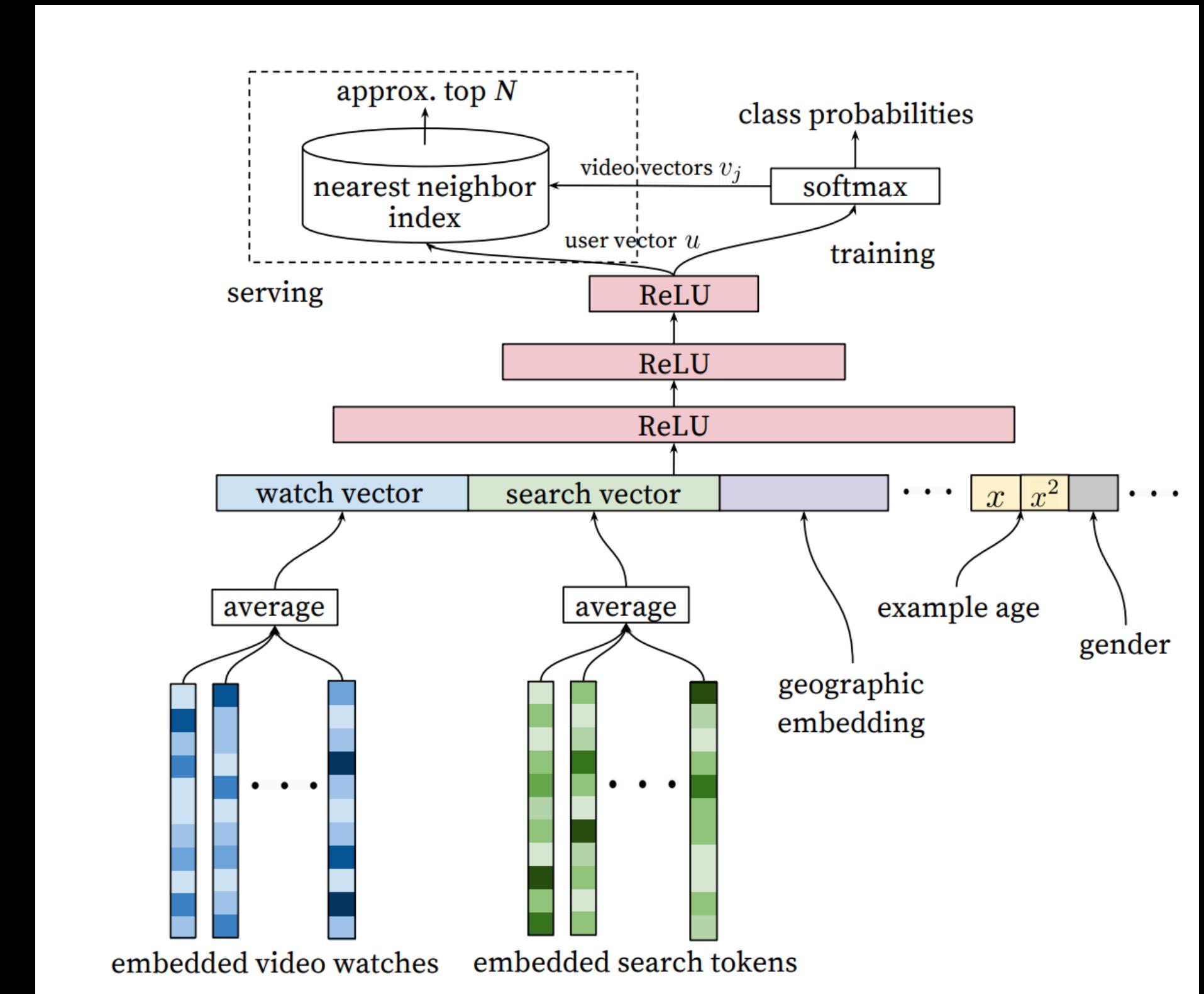
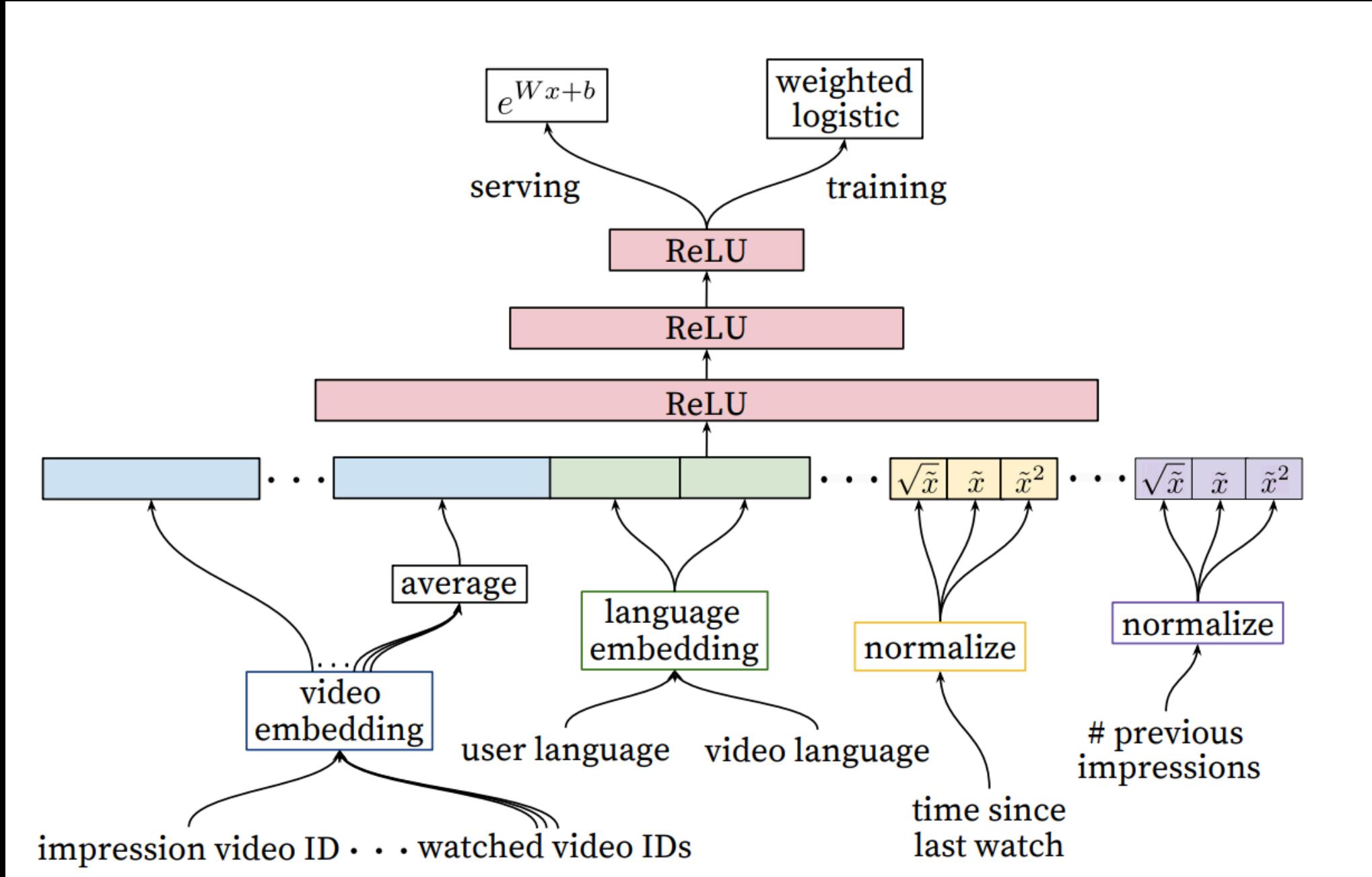
- Better with additional sources then recommended
- ANN in serve time

# YT(youtube) recommender:

- Tips and tricks:
  - Continuous normalisation
  - Shared embedding (last video ID, current video ID, searched video ID)
  - Weighted log regr. For watch time:  
Positive weighted to  $T_i$  watch time  
Negatives unit  
$$\text{odds} \sim \frac{\sum T_i}{N - k},$$
  
learned odds  $\sim E[T](1 + P)$   
(P is small)

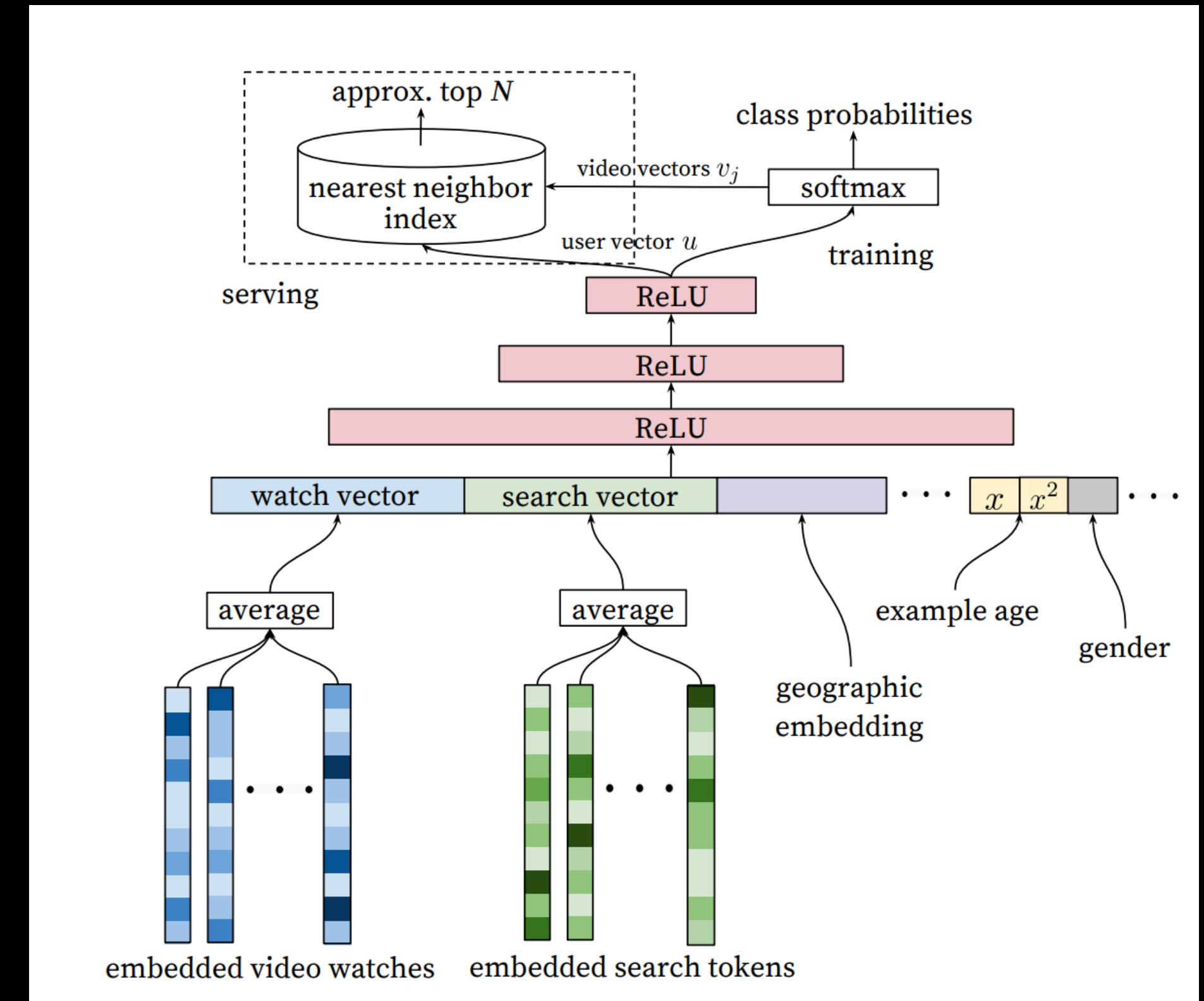
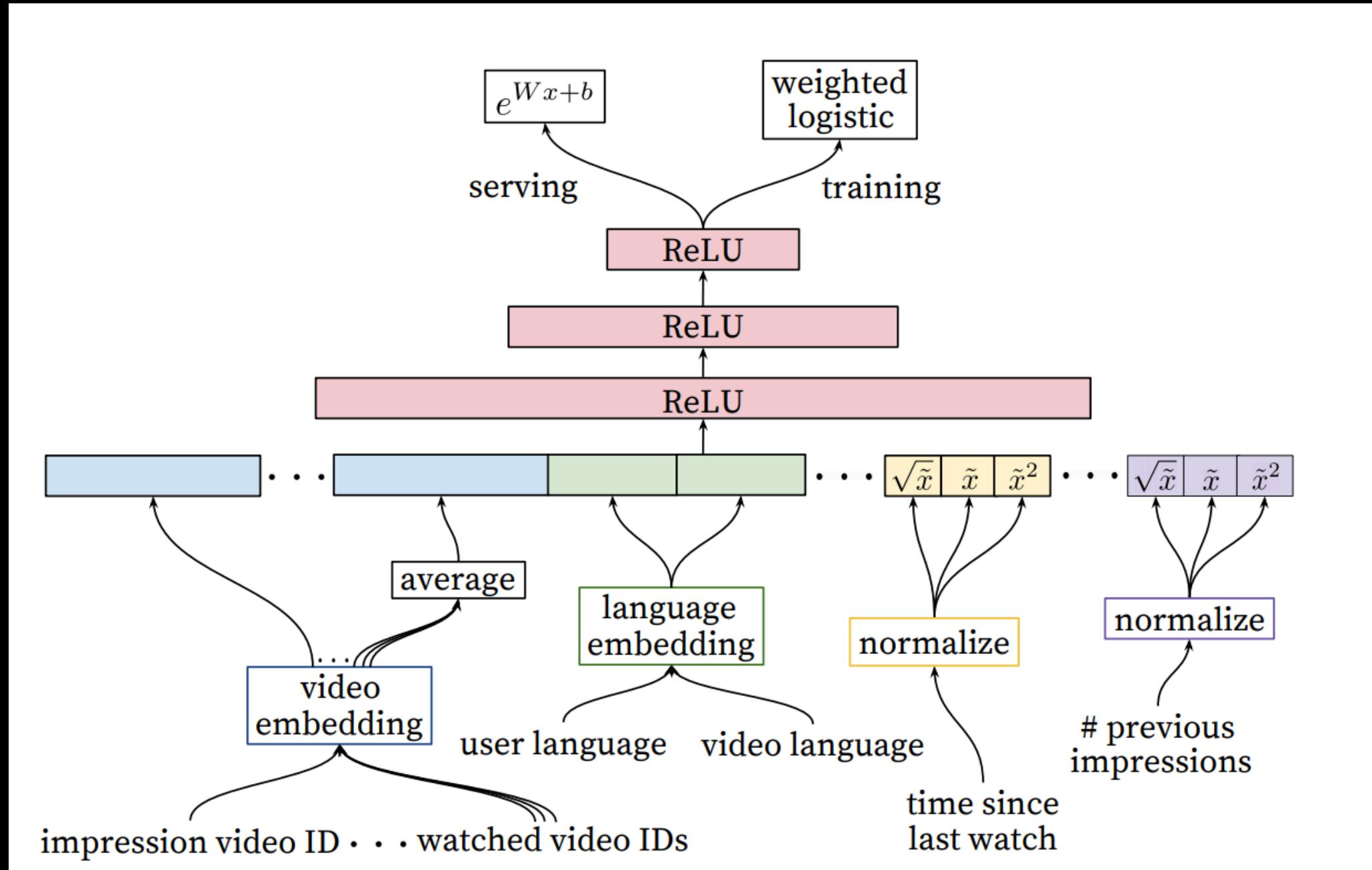


# YT(youtube) recommender:



Problems?

# YT(youtube) recommender:



## Problems?

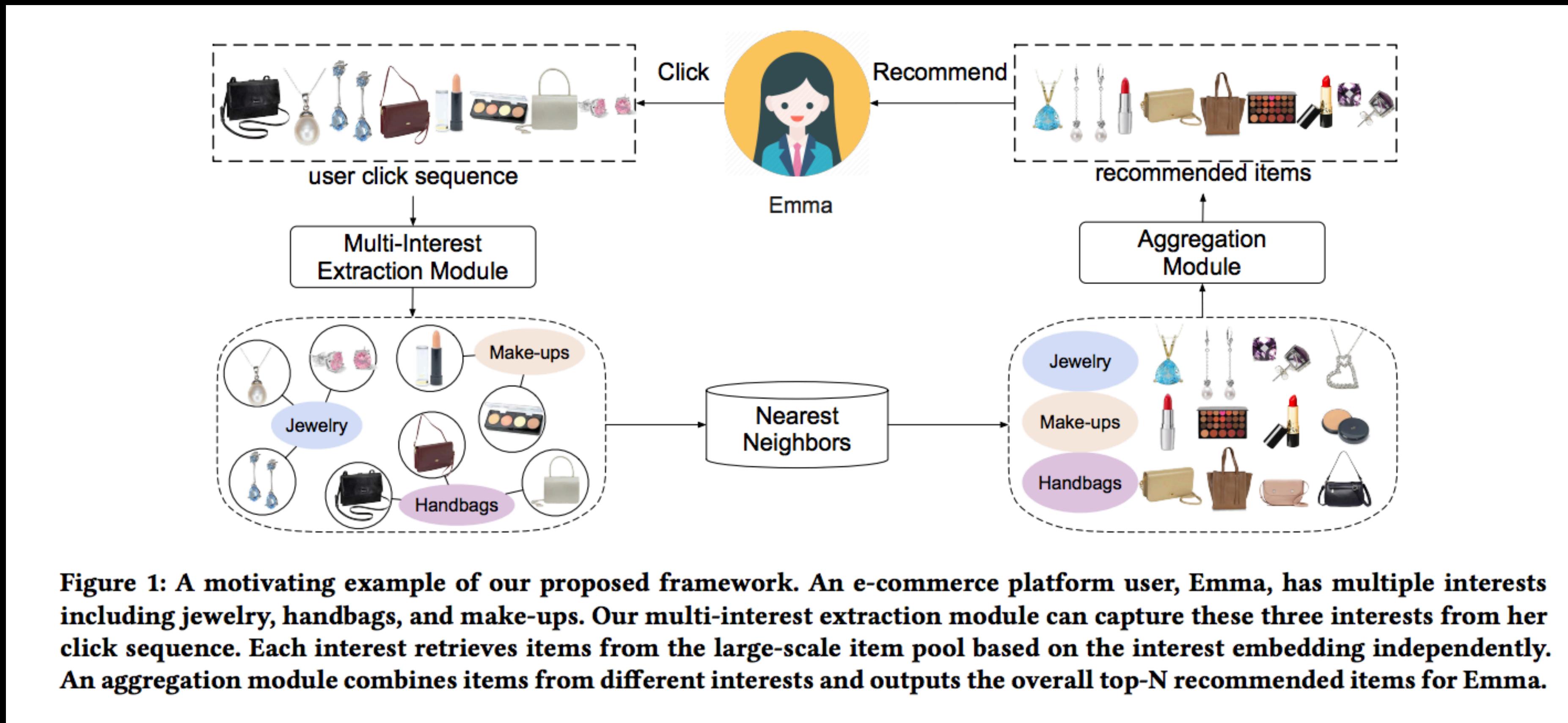
- New items (not users)
- New feedback
- New features
- Time-dependent features

# Recurrent neural networks:

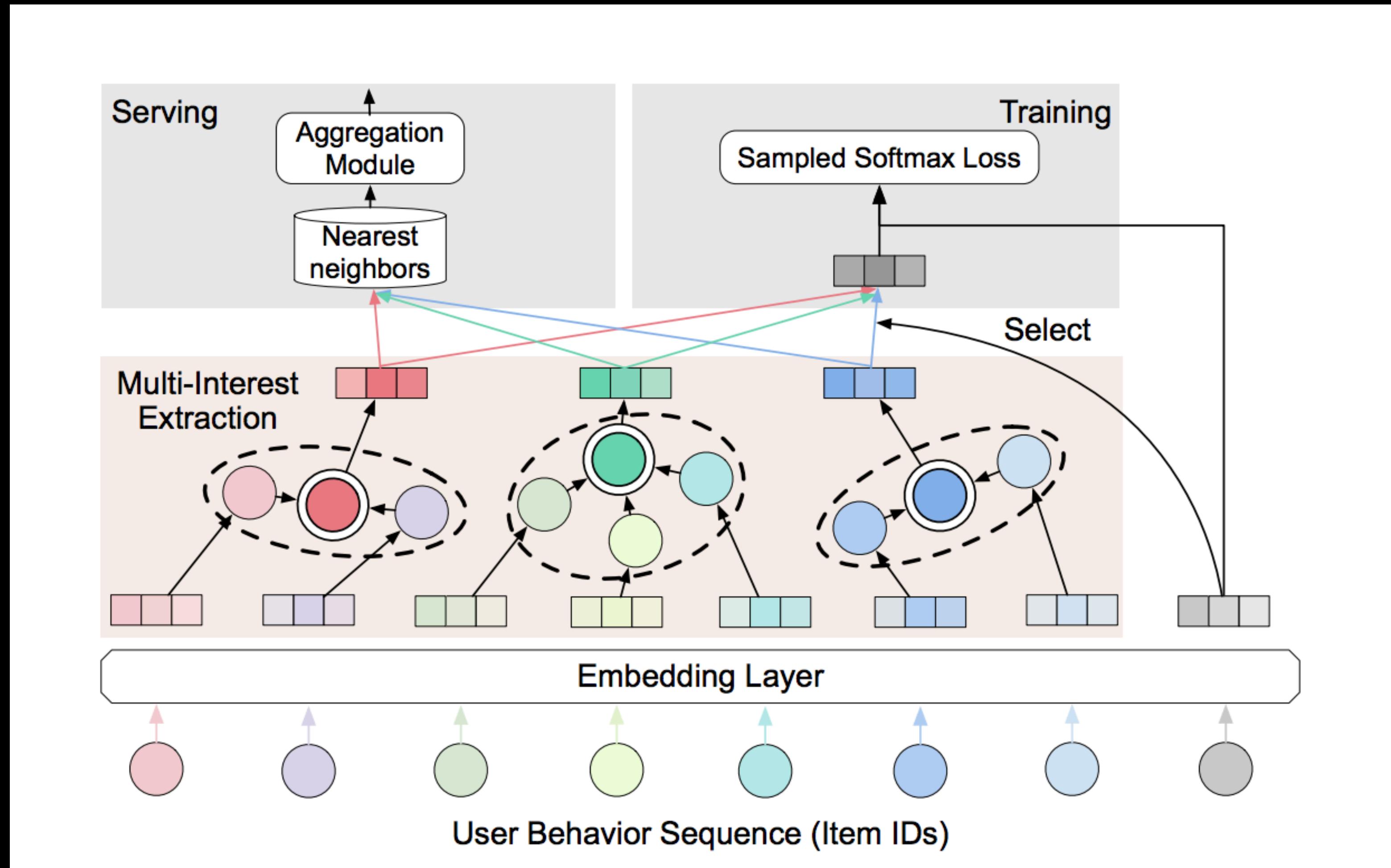


**Ну не заводятся и не заводятся**

# Controllable Multi-Interest Framework for Recommendation:

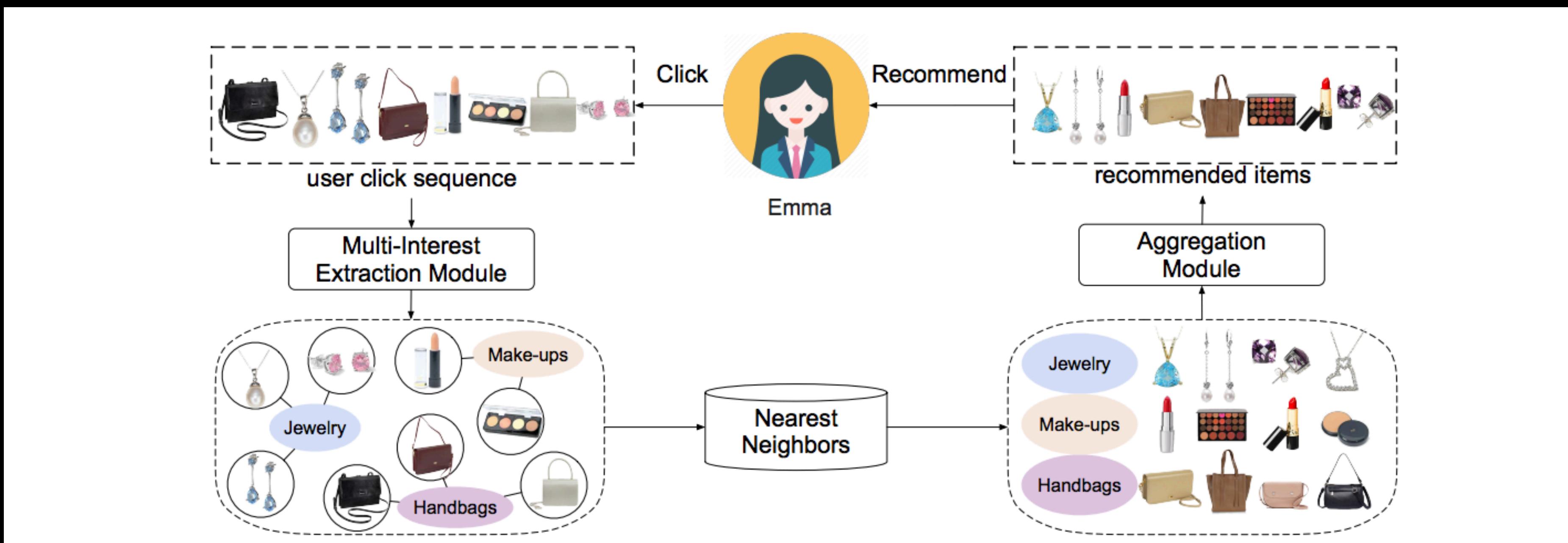


# Controllable Multi-Interest Framework for Recommendation:



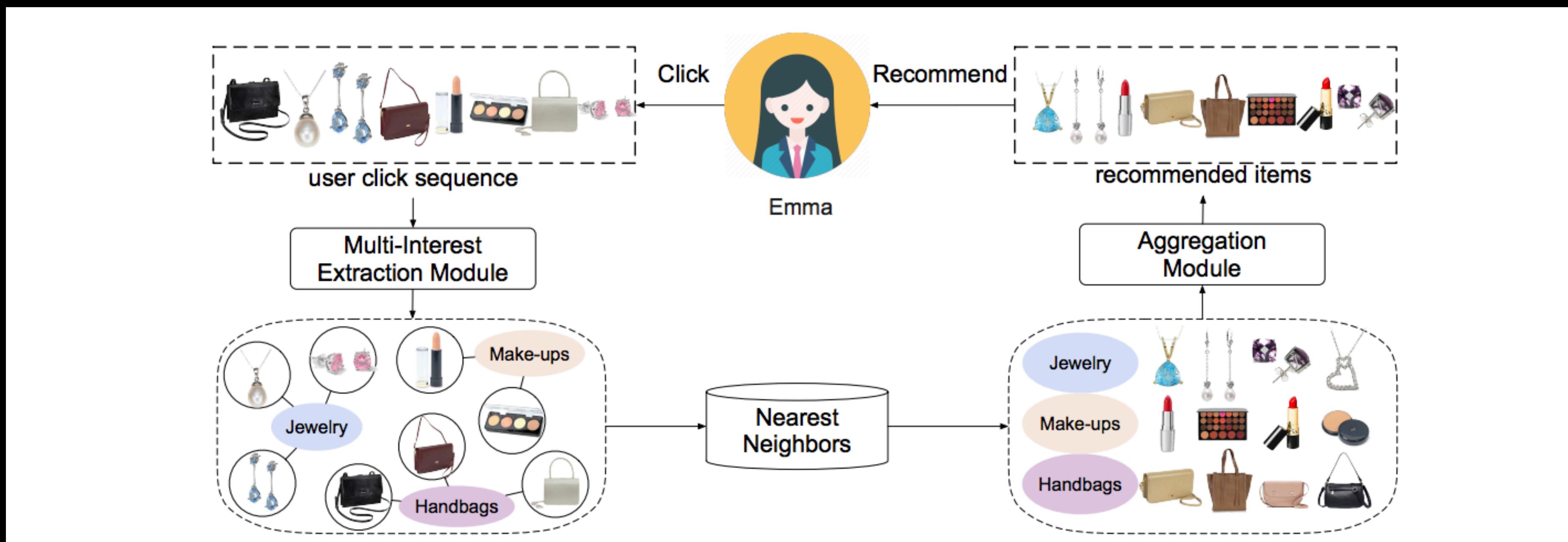
# Controllable Multi-Interest Framework for Recommendation:

- Attention:
  - $A = softmax(W_2^T \tanh(W_1 H^T))$ , where  $H^T$  - prev embeddings matrix,  $H = [d \times n]$ ,  $A = n \times att\_heads$   
 $d$  – embedding size  
 $n$  – history size  
 $att\_heads$  – num heads
  - $V_U = HA$ ;  
 $V_u$  – user interests matrix,  $d \times att\_heads$



# Controllable Multi-Interest Framework for Recommendation:

- Attention:
  - $A = softmax(W_2^T \tanh(W_1 H^T))$ , where  $H^T$  - prev embeddings matrix,  $H = [d \times n]$ ,  $A = n \times att\_heads$   
 $d$  — ebb size  
 $n$  — history size  
 $att\_heads$  — num heads
- Model training:
$$v_u = V_u[:, argmax(V_u^T e_i)]$$
Sampled softmax



# Controllable Multi-Interest Framework for Recommendation:

- Attention:
  - $A = softmax(W_2^T \tanh(W_1 H^T))$ , where  $H^T$  - prev embeddings matrix,  $H = [d \times n]$ ,  $A = n \times att\_heads$   
 $d$  — ebb size  
 $n$  — history size  
 $att\_heads$  — num heads
- Model training:
  - $v_u = V_u[:, argmax(V_u^T e_i)]$   
Sampled softmax
- Aggregation model:
$$f(u, i) = \max_k (e_i^T, v_u^k)$$
$$Q(u, S) = \sum_{i \in S} f(u, i) + \lambda \sum_{i \in S} \sum_{j \in S} g(i, j)$$
$$g(i, j) = 1 \text{ if } CATE(i) \neq CATE(j)$$

# Controllable Multi-Interest Framework for Recommendation:

- Attention:
  - $A = \text{softmax}(W_2^T \tanh(W_1 H^T))$ , where  $H^T$  - prev embeddings matrix,  $H = [d \times n]$ ,  $A = n \times \text{att\_heads}$   
 $d$  — ebb size  
 $n$  — history size  
 $\text{att\_heads}$  — num heads
- Aggregation model:
 
$$f(u, i) = \max_k (e_i^T, v_u^k)$$

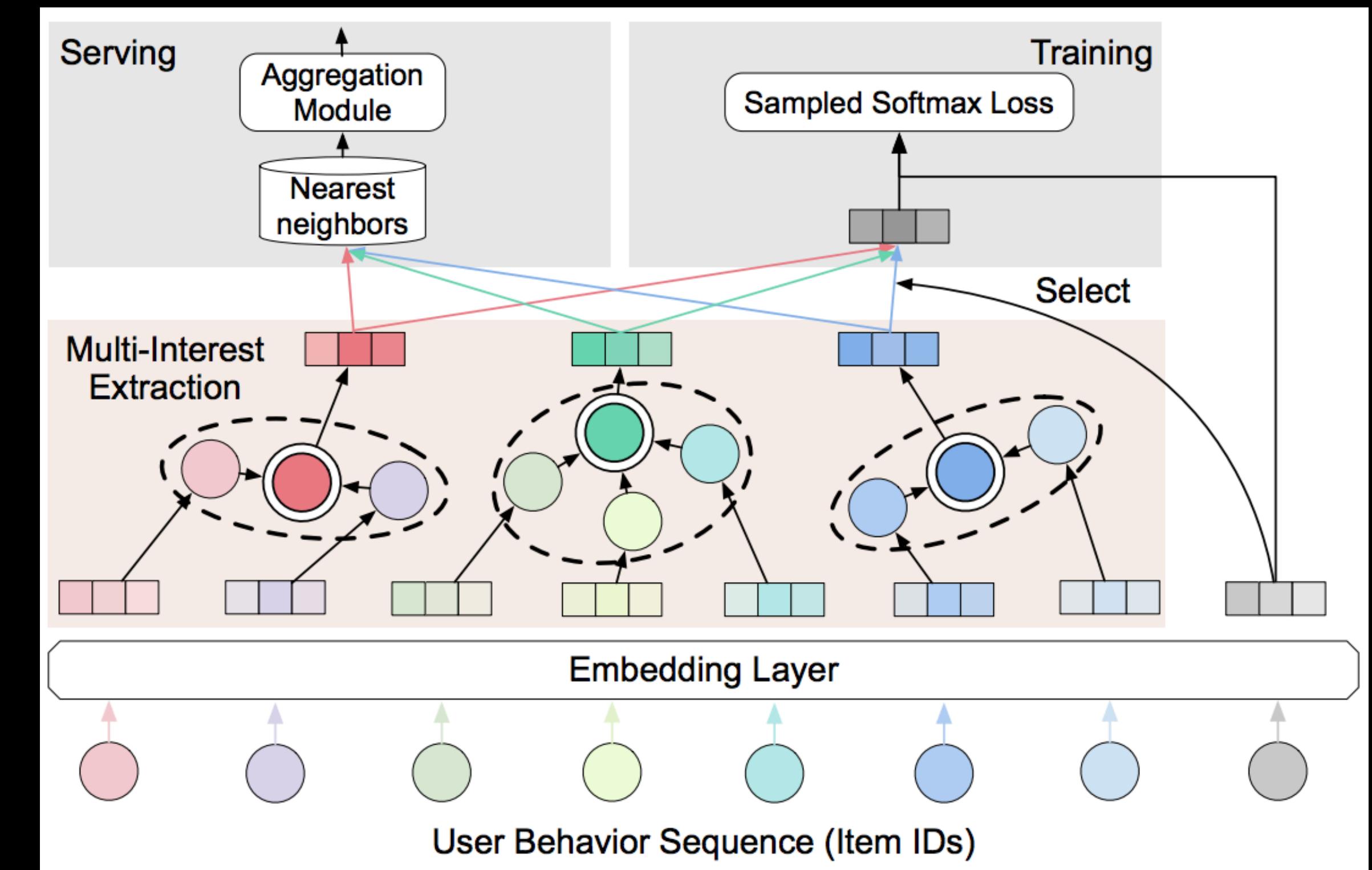
$$Q(u, S) = \sum_{i \in S} f(u, i) + \lambda \sum_{i \in S} \sum_{j \in S} g(i, j)$$

$$g(i, j) = 1 \text{ if } \text{CATE}(i) \neq \text{CATE}(j)$$

**PRBLMS?**

- Model training:
 
$$v_u = V_u[:, \text{argmax}(V_u^T e_i)]$$

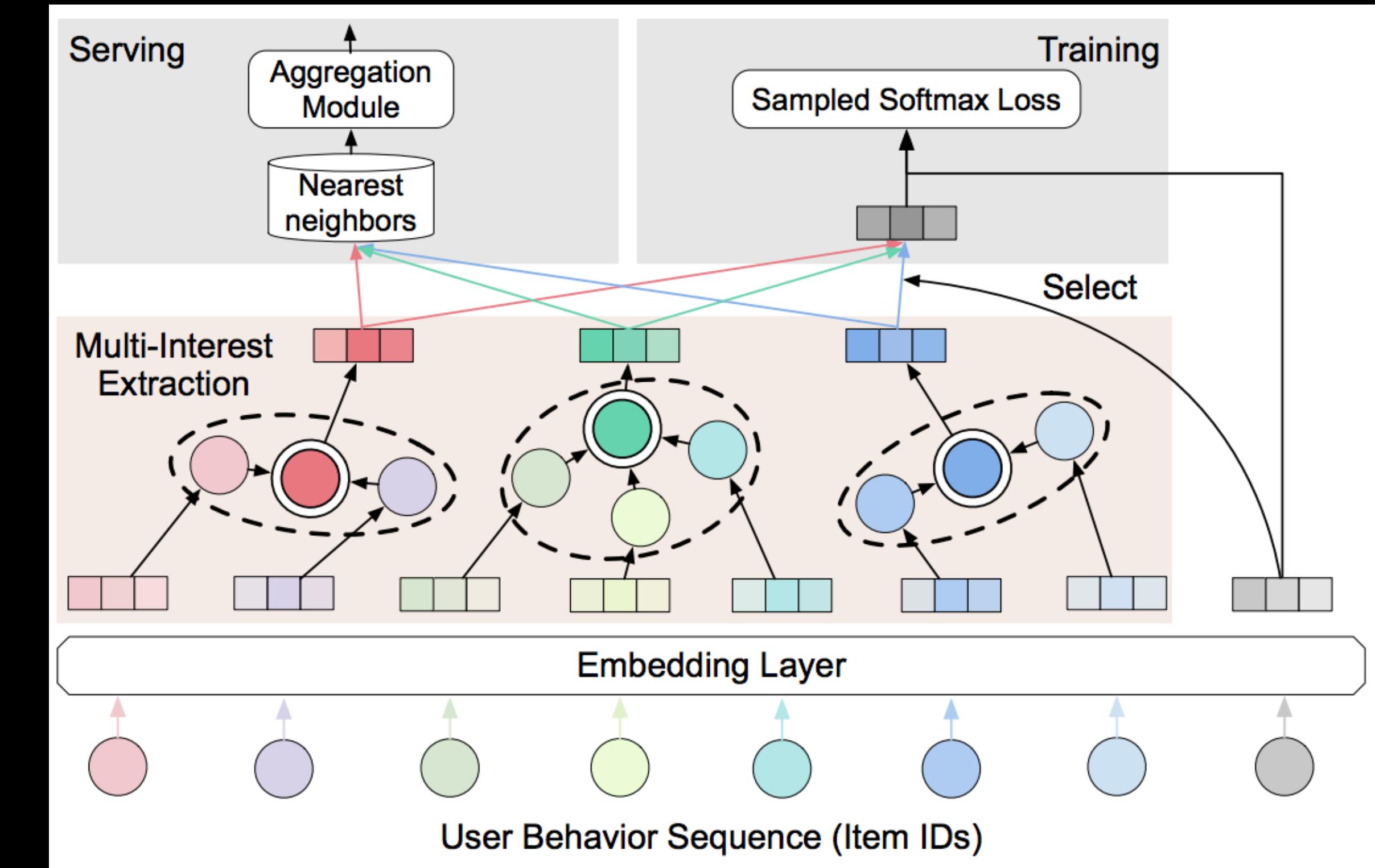
Sampled softmax



# Controllable Multi-Interest Framework for Recommendation:

**PRBLMS?**

- Model training:  
 $v_u = V_u[:, argmax(V_u^T e_i)]$  – same heads
- Pre-setted num of interests
- Over-confident heads
- No feature used



# Controllable Multi-Interest Framework for Recommendation:

PRBLMS?

- Model training:

$$v_u = V_u[:, \text{argmax}(V_u^T e_i)] - \text{same heads}$$

Additive component to loss:

$$\text{loss} = \text{loss} + \alpha * (1 - \text{max\_pair\_dot})$$

- Pre-setted num of interests

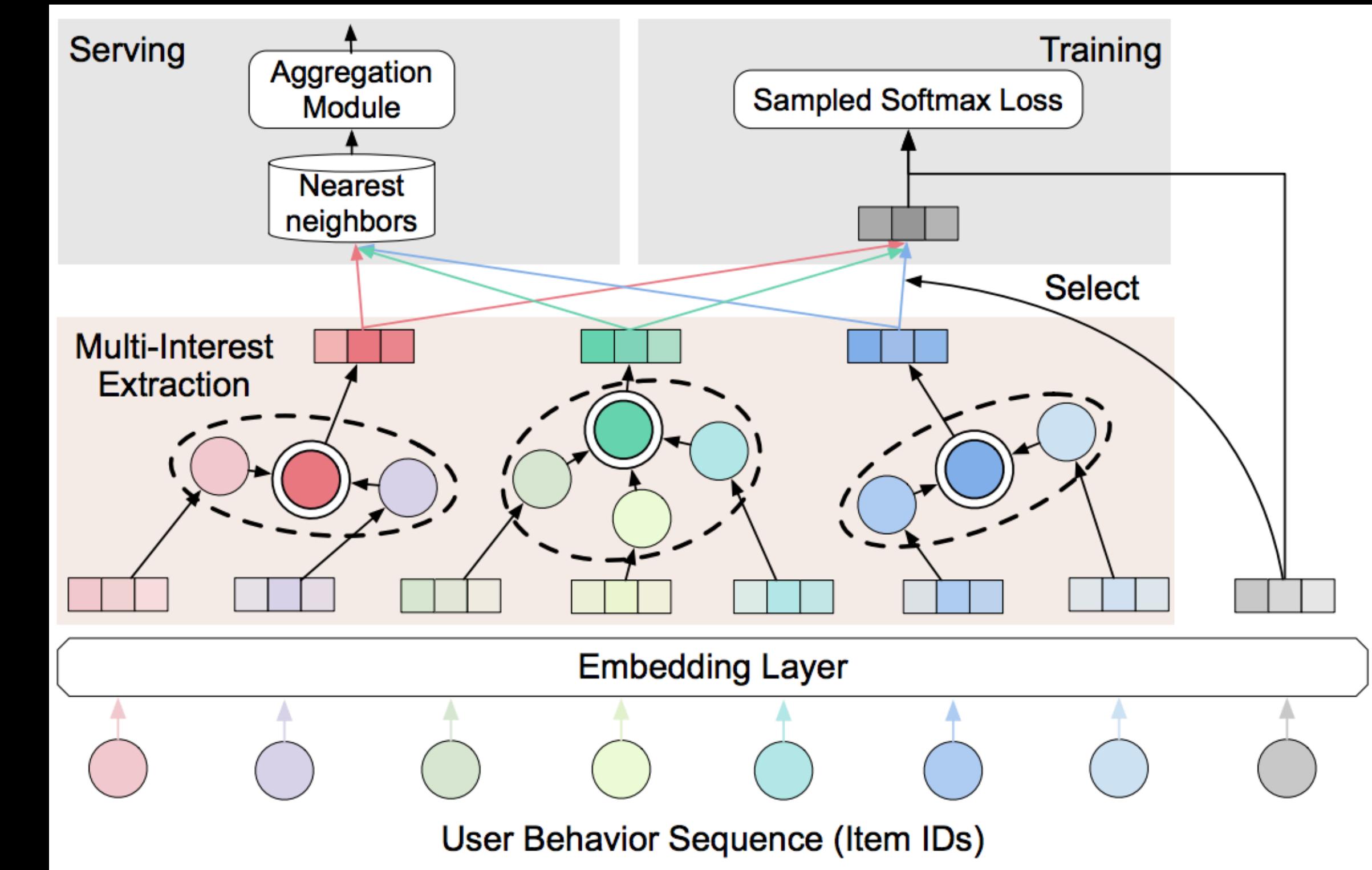
...

- Over-confident heads

Temperature, label smooth, dropout

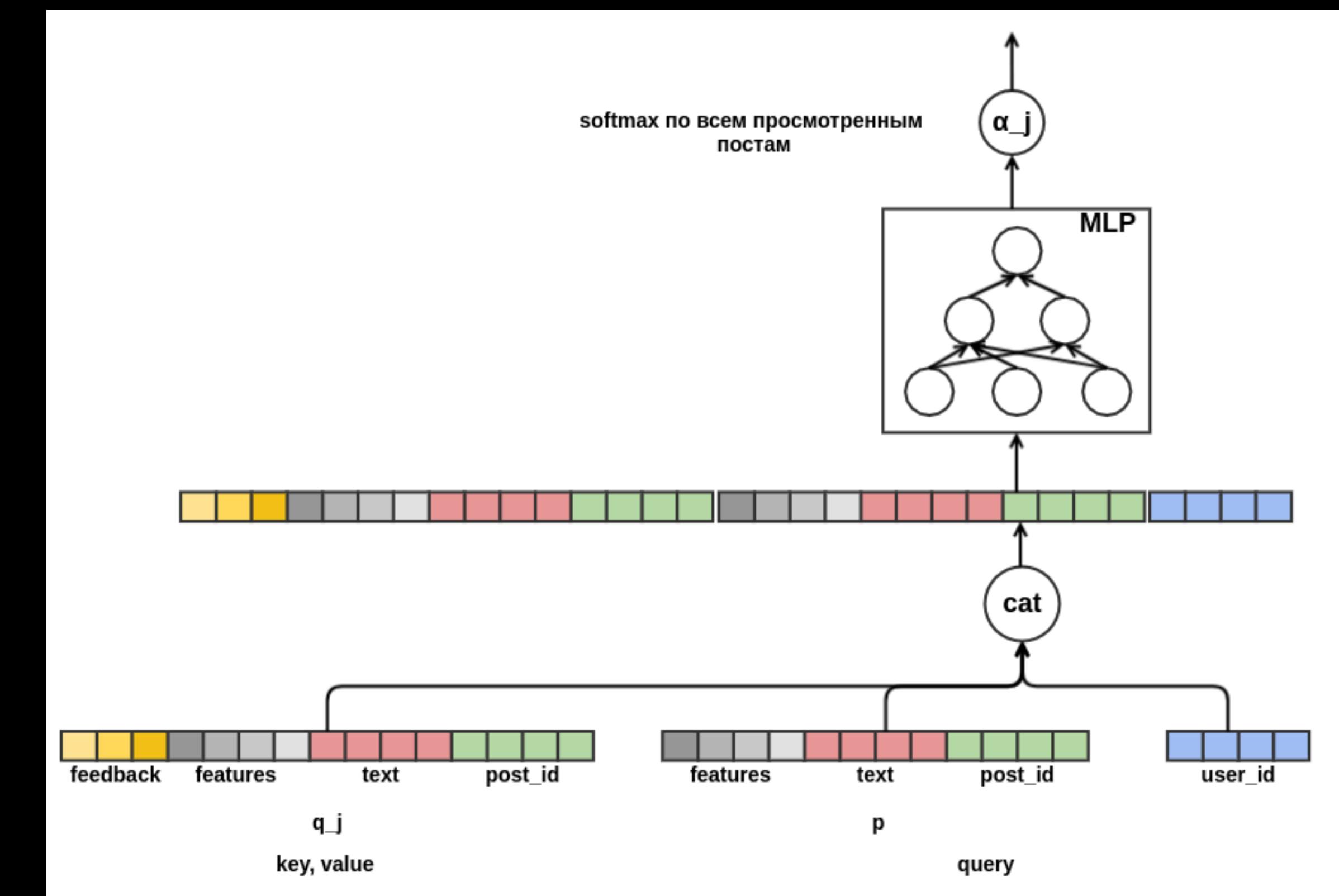
- No feature used

Add some features, it's attention ;)



# Controllable Multi-Interest Framework for Recommendation:

- **No feature used**  
Add some features, it's attention ;)



# MEANTIME: Mixture of Attention Mechanisms with Multi-temporal Embeddings for Sequential Recommendation:

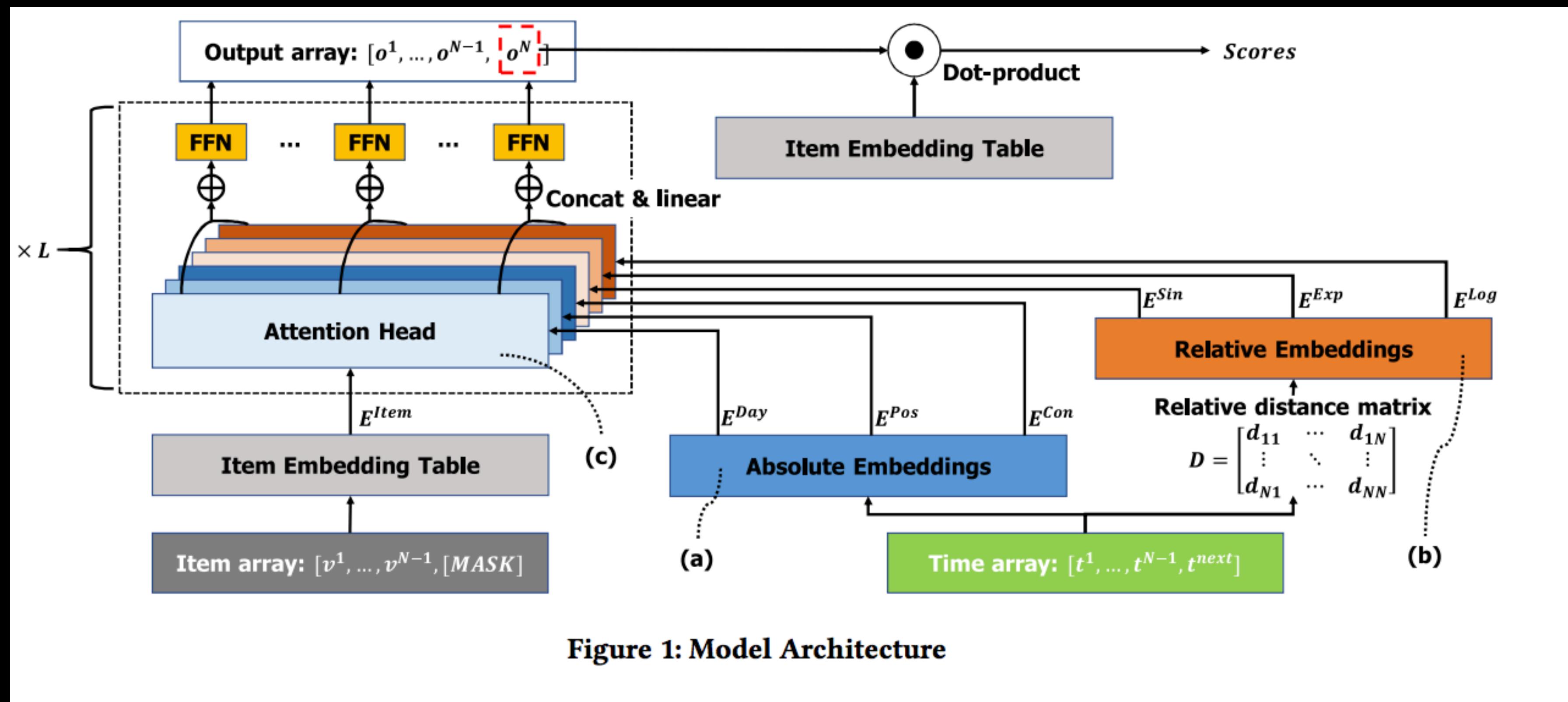


Figure 1: Model Architecture

# MEANTIME: Mixture of Attention Mechanisms with Multi-temporal Embeddings for Sequential Recommendation:

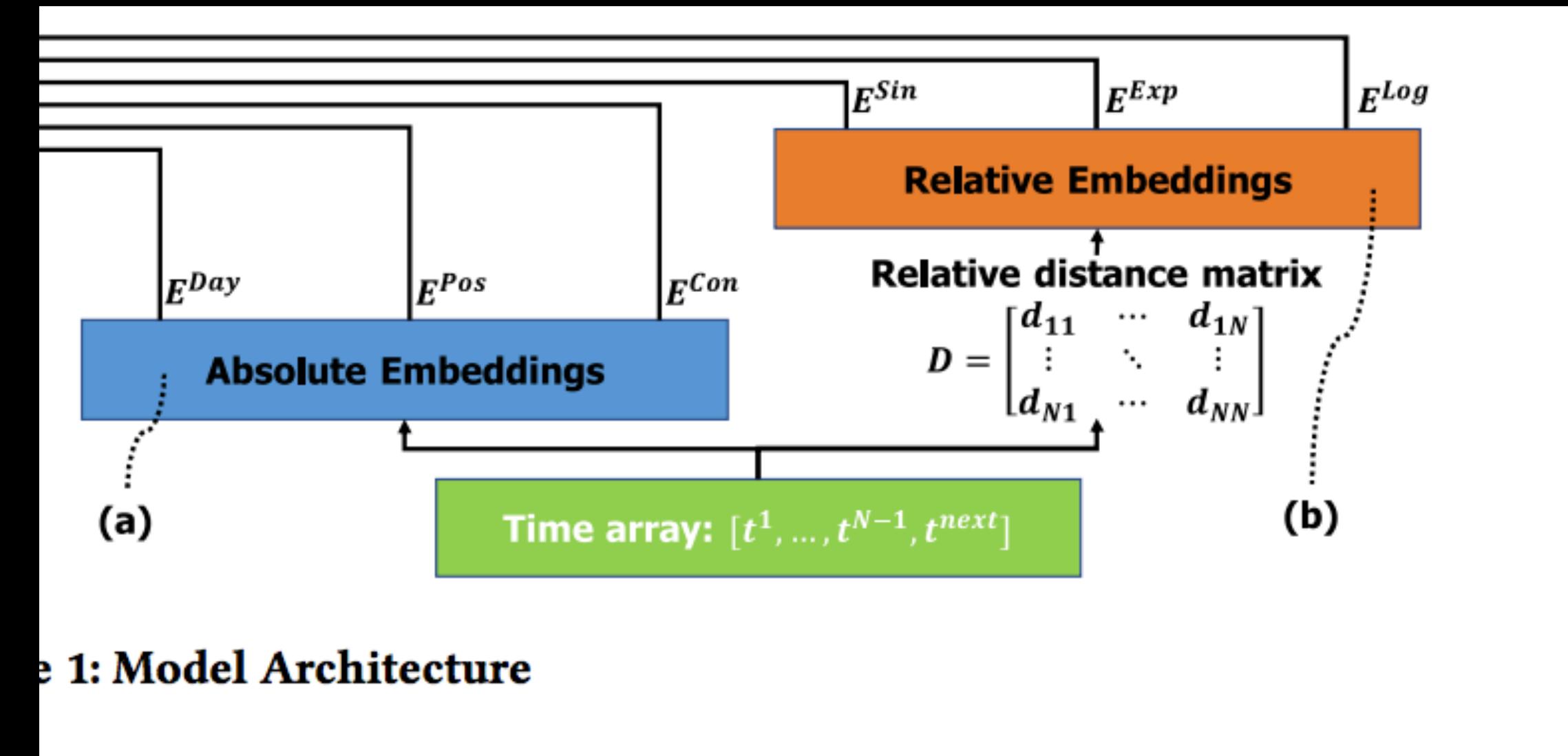
- $E^{pos}$  – relative positional embedding.  
To learn sequential patterns.
- $E^{day}$  – embedding of a day (problems? )  
confine the attention scope within the same  
(or similar) day
- $E^{con}$  – stack of repeated vectors (?) to  
eliminate positional bias

following equation.

$$\vec{\theta}_{ab,2c} = \sin\left(\frac{d_{ab}}{freq^{\frac{2c}{h}}}\right) \quad \vec{\theta}_{ab,2c+1} = \cos\left(\frac{d_{ab}}{freq^{\frac{2c}{h}}}\right) \quad (1)$$

where  $\vec{\theta}_{ab,c}$  is the  $c^{th}$  value of the vector  $\vec{\theta}_{ab}$  and  $freq$  is an adjustable parameter. Likewise, **Exp** encoder and **Log** encoder converts  $d_{ab}$  to  $\vec{e}_{ab}$  and  $\vec{l}_{ab}$  respectively by applying the following equations:

$$\vec{e}_{ab,c} = \exp\left(\frac{-|d_{ab}|}{freq^{\frac{c}{h}}}\right) \quad \vec{l}_{ab,c} = \log\left(1 + \frac{|d_{ab}|}{freq^{\frac{c}{h}}}\right) \quad (2)$$



- Sin captures periodic occurrences.
- Larger time gap is either quickly decayed to zero in **Exp**
- or manageably increased in **Log**.

# Conclusions:

- NN for rank:
  - Accurate feature analysis.
  - Good generalization but lack of interpretability.
  - Not architecture but data analysis.
- It's all about embeddings, but embeddings overfits easily
- Not user embedding, but user embedding as function of items.
- Bag of tricks to make things works: label smooth, time encoders, additives losses and etc.