# Interpretable Click-Through Rate Prediction through Hierarchical Attention

Zeyu Li, Wei Cheng, Yang Chen, Haicheng Chen, Wei Wang

UCLA, NEC Labs America, Google

WSDM 2020

Houston, Texas, USA

# Recommender systems

# Click-throughs

Two natural questions to ask:

1. How many advertisements will be clicked?
2. How many clicks will be purchased?

Two natural questions to ask:

1. How many advertisements will be clicked?
2. How many clicks will be purchased?

o CTR:

   o Important role in recommendation system

   o Revenue of advertisements



**Image**: https://www.lyfemarketing.com/blog/average-click-through-rate/

# Background

o CTR: *binary prediction*

o Pre-Deep Learning Model
- o FM: Factorization Machine
- o MF: Matrix Factorization
- o LR: Logistic Regression

o Deep learning based CTR model
- o DeepFM  = FM module  + Deep module
- o xDeepFM = CIN module + Deep module
  - o CIN: Compressed Interest Network
- o and more …

**Figure 5: The architecture of xDeepFM.**

Figure 1: Wide & deep architecture of DeepFM. The wide and deep component share the same input raw feature vector, which enables DeepFM to learn low- and high-order feature interactions simultaneously from the input raw features.

| xDeepFM | DeepFM |
|---|---|
| WDN | DCN | PNN |



**Wide & Deep Models**

**Figure 1: The Deep & Cross Network**

Product-based Neural Network Architecture.

input layer

hidden layer 1    hidden layer 2

output layer

o DNN

 o Widely used in CTR models

 o Unjustifiable element-wise computation within representations of input or hidden features

 o Unaffordable complexity for big feature dim or size

**Image**: https://hackernoon.com/challenges-in-deep-learning-57bbf6e73bb

Samueli
**Computer Science**
SCALABLE ANALYTICS INSTITUTE
NEC LABORATORIES AMERICA, INC.
*Relentless* passion for innovation

# Concerns of DNN

o Okay for online shopping with general purposes
   o Shopping on Amazon …

o *NOT* okay for:
   o Medicine recommendation
   o Financial service recommendation

o Criteo:
   o 4 billions clicks in 24 hrs

o Interpretability

o Attention mechanism

o Avoid flat concatenation of features

o Avoid DNN and dim-wise computation

o Efficiency

o Shrunk problem size

o Self attention in Transformer



**Right Figure**: Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems. 2017.

# Hierarchical Attention

○ Input: i-th order features:

○ Generate aggregated feature:

$$\alpha_i^{(j)} = \frac{\exp{(\boldsymbol{c}_i^T ReLU(\mathbf{W}_i \boldsymbol{x}_i^{(j)})}}{\sum_{j' \in F} \exp{(\boldsymbol{c}_i^T ReLU(\mathbf{W}_i \boldsymbol{x}_i^{(j')}))}},$$

$$\boldsymbol{u}_i = \text{AttentionalAgg}(\mathbf{X}_i) = \sum_{j=1}^{m} \alpha_i^{(j)} \boldsymbol{x}_i^{(j)},$$

○ Output (i+1)-t order features:

$$\boldsymbol{x}_{i+1}^{(j)} = \boldsymbol{u}_i \circ \boldsymbol{x}_1^{(j)} + \boldsymbol{x}_i^{(j)}, \quad j \in \{1, \ldots, m\},$$

Samueli
**Computer Science**
ScAi SCALABLE ANALYTICS INSTITUTE
NEC
NEC LABORATORIES AMERICA, INC.
*Relentless* passion for innovation
13

o # Datasets

  o ## Performance evaluation

| Dataset | Criteo | Avazu | Frappe |
|---|---|---|---|
| #. of features (C + N) | 22 + 14 | 21 + 0 | 7 + 0 |
| #. of total records | 13.8M | 12.1M | 288K |
| #. of distinct features | 605.7K | 23.8K | 5,382 |

    o Critio, Avazu, Frappe

  o ## Interpretability study

    o Movielens-1m dataset (reviews as clicks)

o # Baselines

  o # FM, Wide&Deep, DCN, PNN, DeepFM, xDeepFM

o # Metrics

  o # Area Under ROC Curve (AUC)
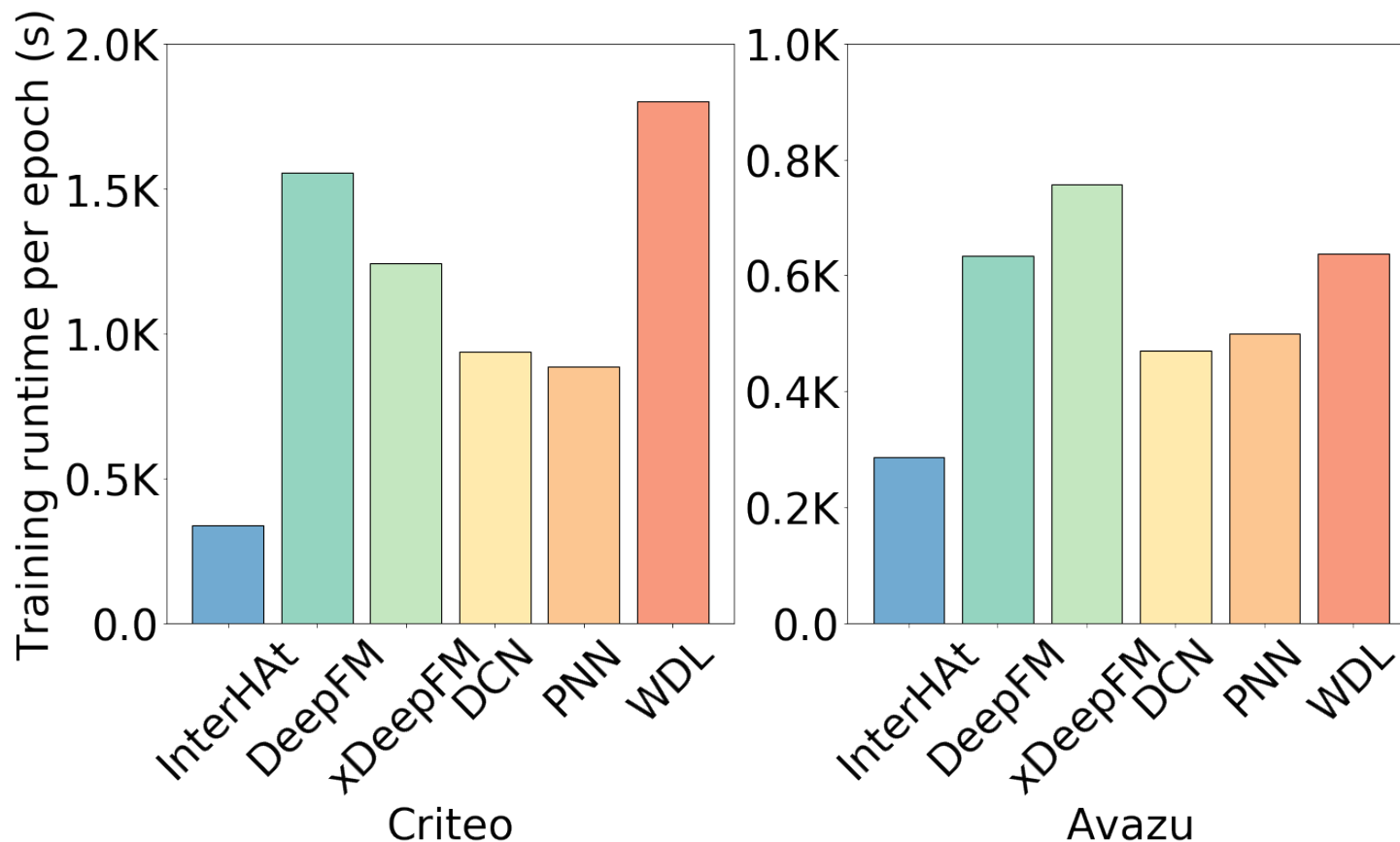
  o # Cross Entropy (LogLoss)

# Performance

o Comparable with SOTA models

o Perform better on categorical features
   - o SOTA models have close performance
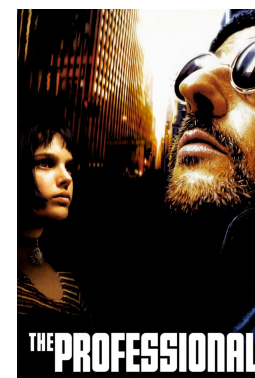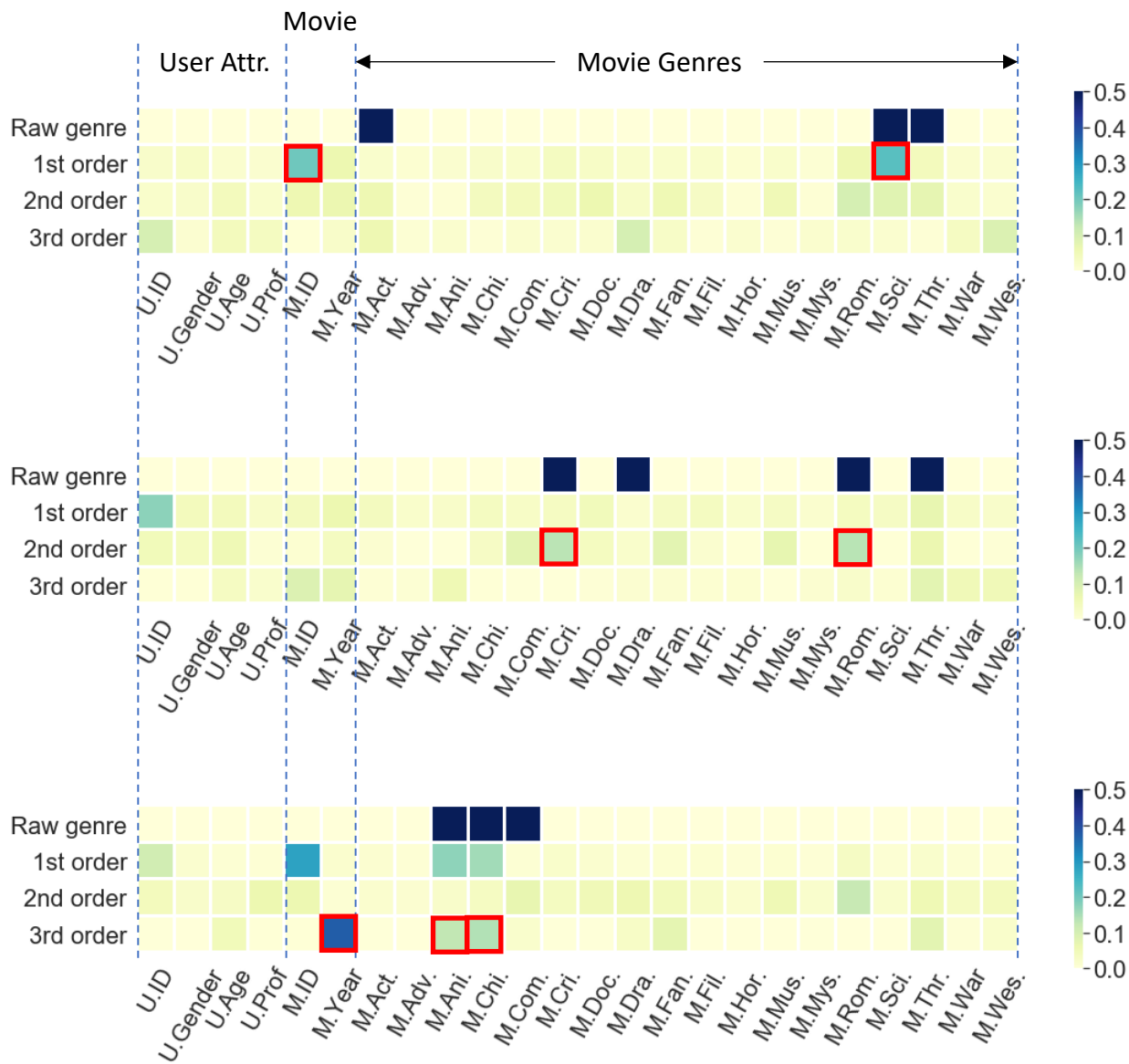   - o Need better ways for encoding numeric features

| Dataset | Criteo | | Avazu | | Frappe | |
|---|---|---|---|---|---|---|
| Metrics | Logloss | AUC | Logloss | AUC | Logloss | AUC |
| FM | 0.4814 | 0.7525 | 0.3951 | 0.7508 | 0.4480 | 0.8625 |
| Wide&Deep | 0.4577 | 0.7845 | 0.3920 | 0.7564 | 0.2571 | 0.9500 |
| DCN | 0.4590 | 0.7826 | 0.3921 | 0.7564 | 0.2335 | 0.9616 |
| PNN | **0.4547** | **0.7887** | 0.3916 | 0.7569 | 0.2177 | 0.9642 |
| DeepFM | 0.4560 | 0.7866 | 0.3920 | 0.7561 | 0.2410 | 0.9520 |
| xDeepFM | 0.4563 | 0.7874 | 0.3917 | 0.7569 | 0.2043 | 0.9694 |
| InterHAt-S | 0.4608 | 0.7820 | 0.3919 | 0.7577 | 0.2151 | 0.9616 |
| InterHAt | 0.4577 | 0.7845 | **0.3910** | **0.7582** | **0.2026** | **0.9696** |

o InterHAt trains faster than other baselines

# Interpretability

# Conclusion

o InterHAt:

  o Efficiency and interpretability issues of CTR task

    o Efficiency:

      o Avoiding **deep** fully connect neural networks

    o Interpretability:

      o Attention mechanism

      o Interpretability v.s. Explanability

  o Nice performances on both aspects!

  o Try it out:

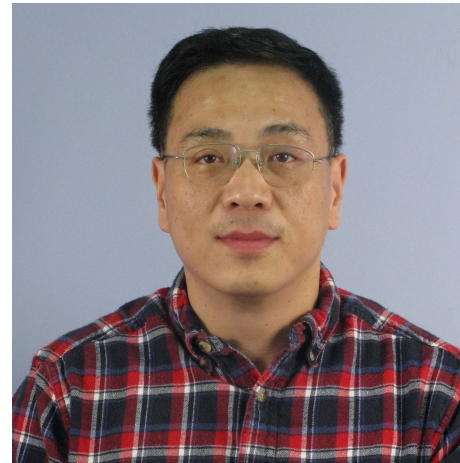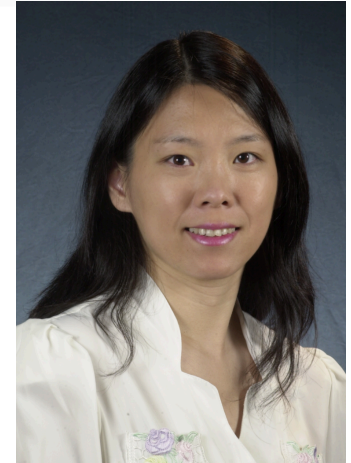    o `https://github.com/zyli93/InterHAt`

# Questions?



Wei Cheng          Yang Chen          Haifeng Chen          Wei Wang