

# Hierarchical Context enabled Recurrent Neural Network for Recommendation

Kyungwoo Song\*, Mingi Ji\*, Sungrae Park, Il-Chul Moon

(\* : Equal Contribution)



## 30-Second Summary

**Question :** Can we detect where the **user's interest changes**?

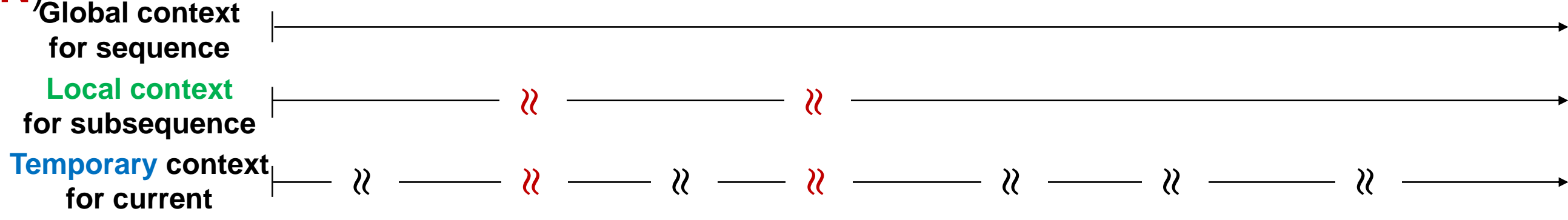
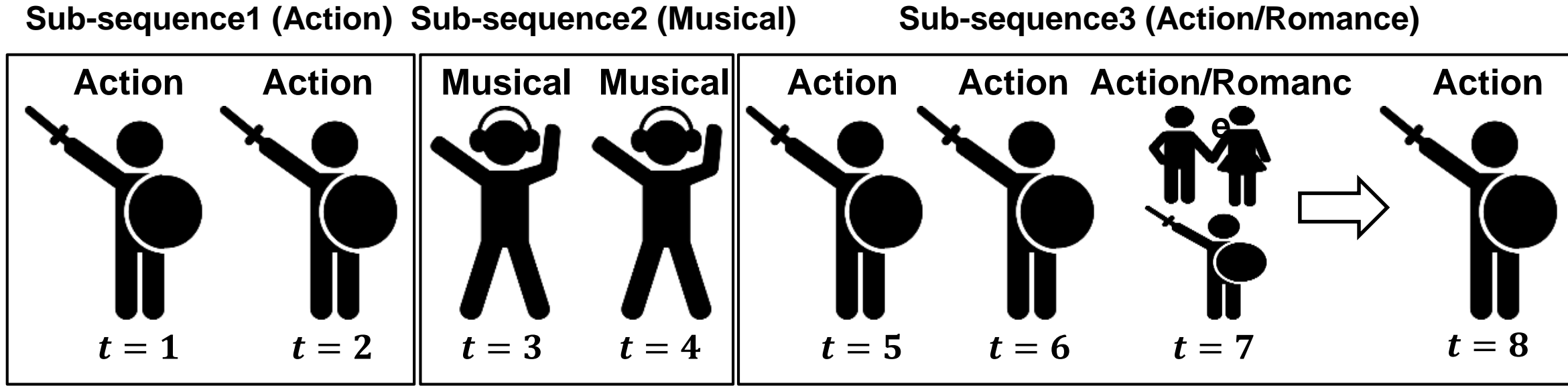
**Our answer :** Yes!

**How? :** **Interest drift assumption**

- “If the user's **local context** (for sub-sequence) and the current item are very different, the user's **temporary** interest drift occurs.”

**More specific :** Hierarchical Context enabled Recurrent Neural Network (**HCRNN**)

- Incorporate the interest drift assumption
- Design **hierarchical contexts** (global, local, and temporary)
- Keep **local** and **temporary** contexts independently
- Introduce **interest drift gate** to capture the interest drift



## Motivation

- A user history is a sequence of user orders or clicks, and the history represents the user's interest
- A long user history inevitably reflects the transitions of personal interests over time
- We can predict next item better if we include modeling on an interest drift of users.

## Model Assumption

- The user's interest can be hierarchically ranging from general interest to a temporary (**global**, **local**, and **temporary**)
- Each hierarchical context have different abstract levels of information.
- Interest drift assumption**

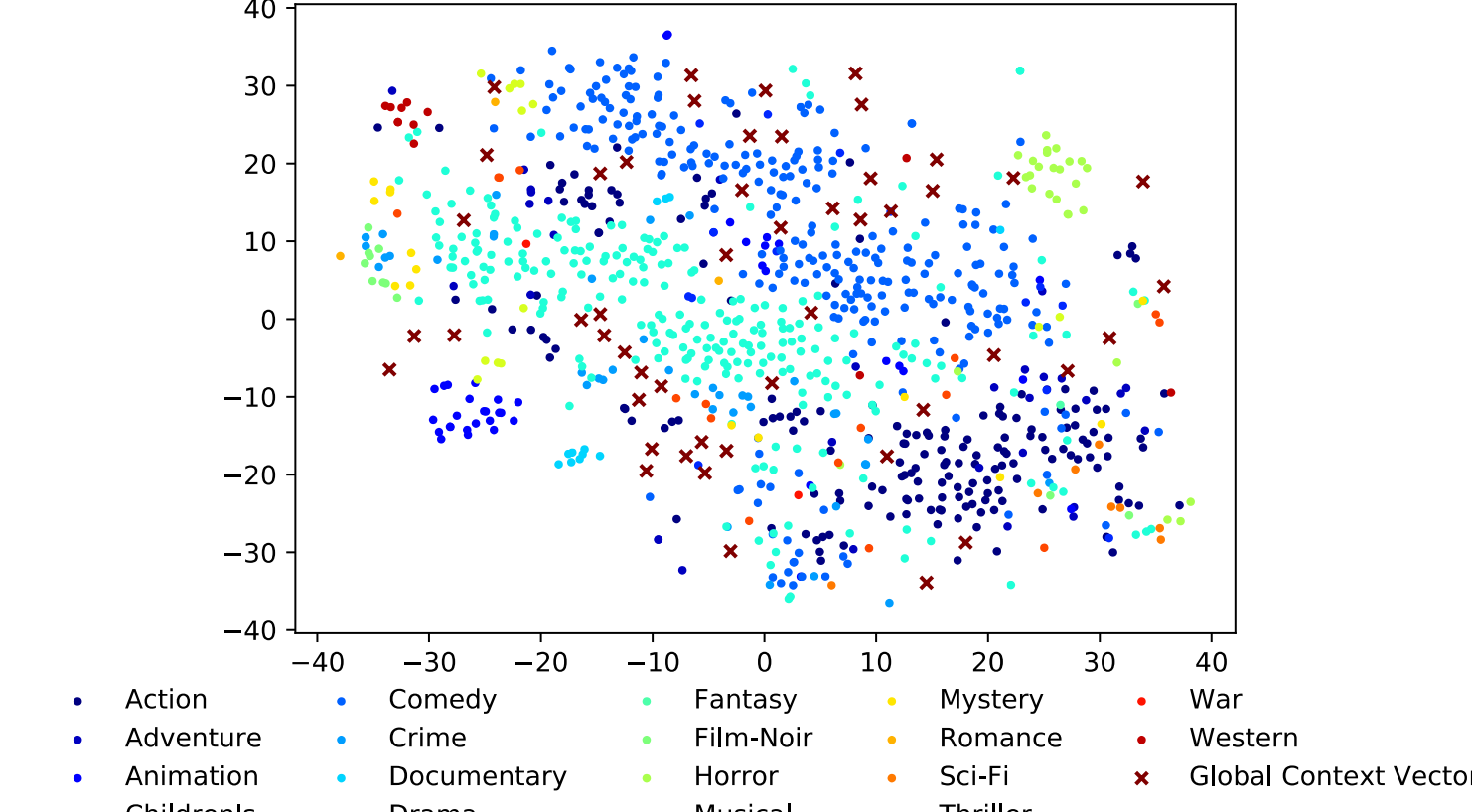
## Results

### 1) Quantitative Results

	CiteULike				LastFM				MovieLens			
	R@3	R@20	M@3	M@20	R@3	R@20	M@3	M@20	R@3	R@20	M@3	M@20
POP	1.44	5.78	0.92	1.44	0.37	1.99	0.34	0.51	2.43	12.51	1.54	2.65
S-POP	1.26	4.99	0.79	1.23	0.87	3.65	0.55	0.87	2.27	12.23	1.42	2.52
Item-KNN	0.00	6.90	0.00	4.79	0.00	11.59	0.00	8.00	0.00	6.32	0.00	4.28
BPR-MF	0.49	3.15	0.27	0.60	0.82	2.15	0.59	0.73	1.69	8.93	1.07	1.91
LSTM4REC	7.07	23.33	4.93	6.82	15.29	24.75	12.68	13.95	8.52	32.80	5.63	8.45
GRU4REC	8.37	24.19	5.98	7.86	18.29	26.46	13.85	16.95	8.50	32.74	5.60	8.42
NARM	7.81	24.82	5.40	7.41	18.30	33.60	13.12	15.25	9.14	33.42	6.09	8.93
STAMP	5.09	21.93	3.25	5.22	9.29	19.84	6.62	8.01	3.95	20.52	2.65	4.47
HCRNN-1	8.60	25.36	6.18	8.16	20.67	34.40	15.77	17.68	9.23	33.78	6.13	9.00
HCRNN-2	8.83	25.10	6.41	8.38	20.78	34.14	16.20	18.08	9.22	33.76	6.14	9.01
HCRNN-3	9.21	25.42	6.65	8.61	21.39	34.72	16.66	18.52	9.38	33.67	6.23	9.08
HCRNN-3+Bi	<b>9.33*</b>	<b>25.81*</b>	<b>6.74*</b>	<b>8.70*</b>	<b>21.90*</b>	<b>34.80*</b>	<b>17.33*</b>	<b>19.12*</b>	<b>9.53*</b>	<b>33.83*</b>	<b>6.38*</b>	<b>9.21*</b>
Improvement (%)	11.47	3.99	12.71	10.69	19.67	3.57	9.34	12.80	4.27	1.23	4.76	3.14

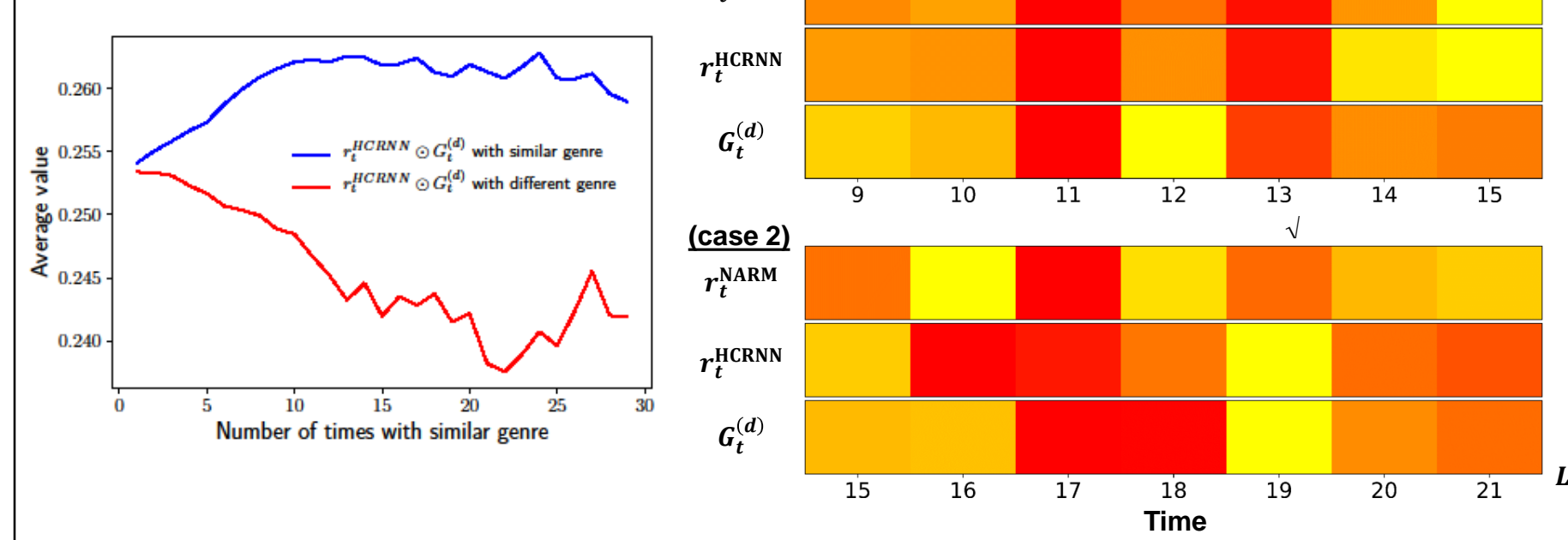
- HCRNN-1 > Baselines (NARM, STAMP)
- necessity of hierarchical context
- HCRNN-3 > HCRNN-2, HCRNN-1
- Interest drift assumption is experimentally justifiable.
- HCRNN-3+Bi > HCRNN-3
- bi-channel attention with hierarchical contexts may improve the performance experimentally.

### 2) Context Embedding



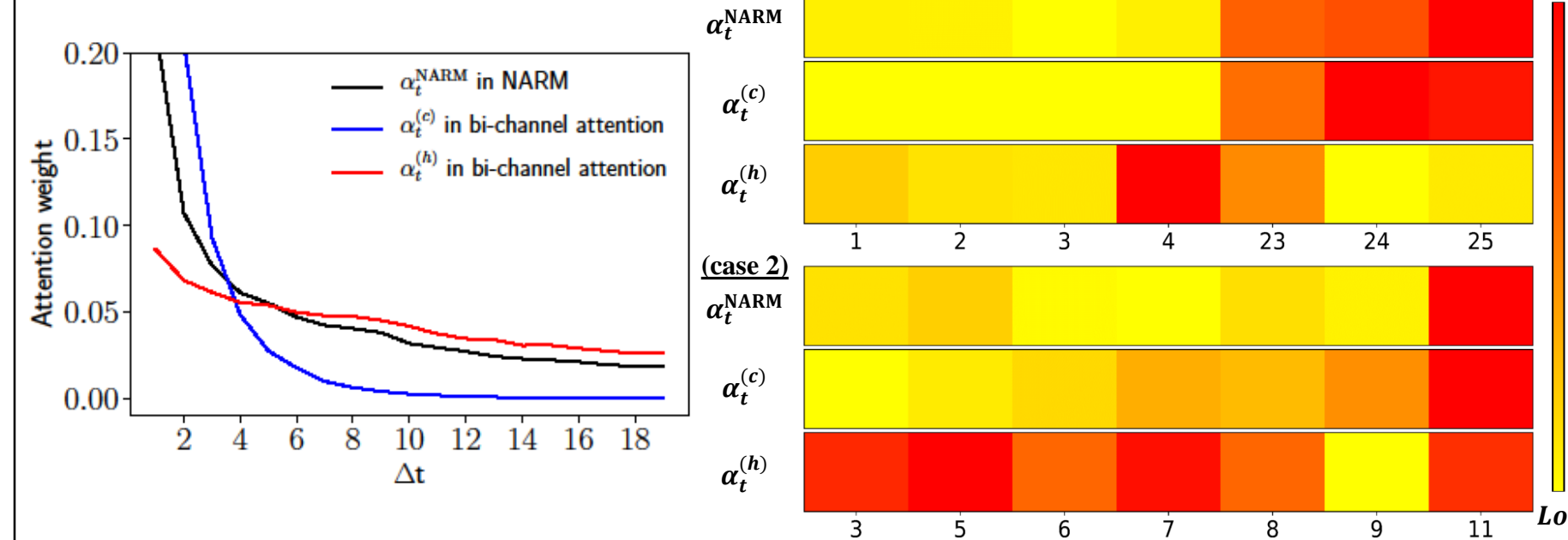
	Genre	Movie Title
$M_{global}^{(6)}$	Animation	Pinocchio, Yellow Submarine, Snow White and the Seven Dwarfs
$M_{global}^{(19)}$	Action	Star Trek: Generations, Predator, Butch Cassidy and the Sundance Kid
$M_{global}^{(31)}$	Horror	Scream, An American Werewolf in London, Dracula

### 3) Gate Analysis



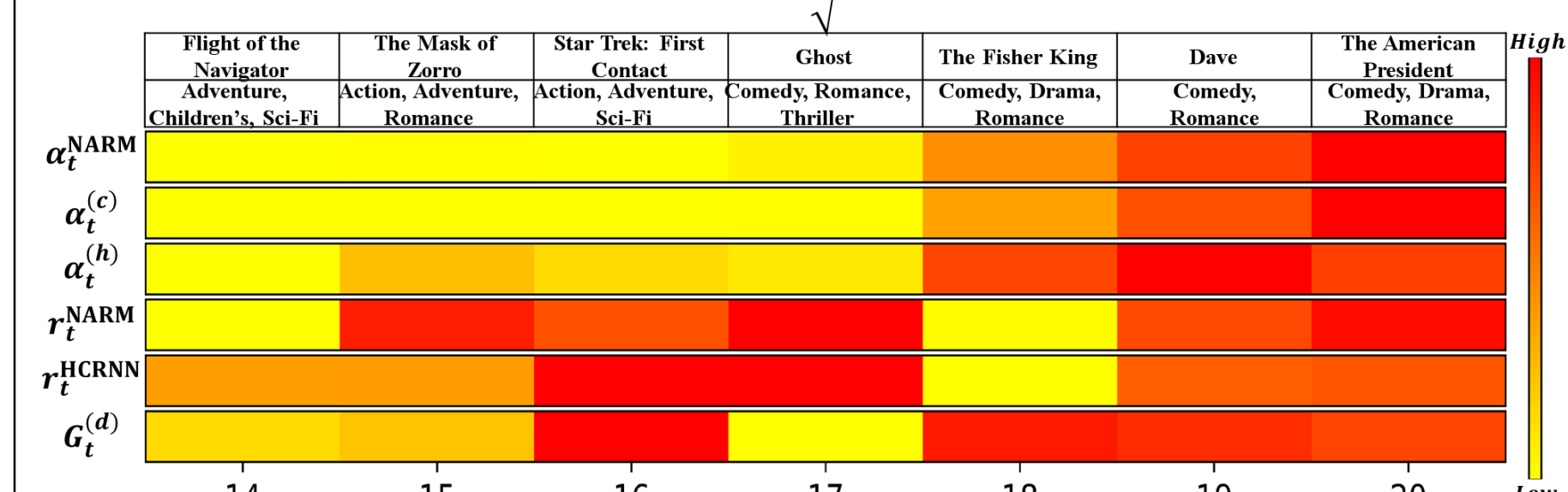
- If the genre of the current input is different with previous items,  $r_t \odot G_t^{(d)}$  has a smaller value compared to the opposite situation.

### 4) Bi-Channel Attention



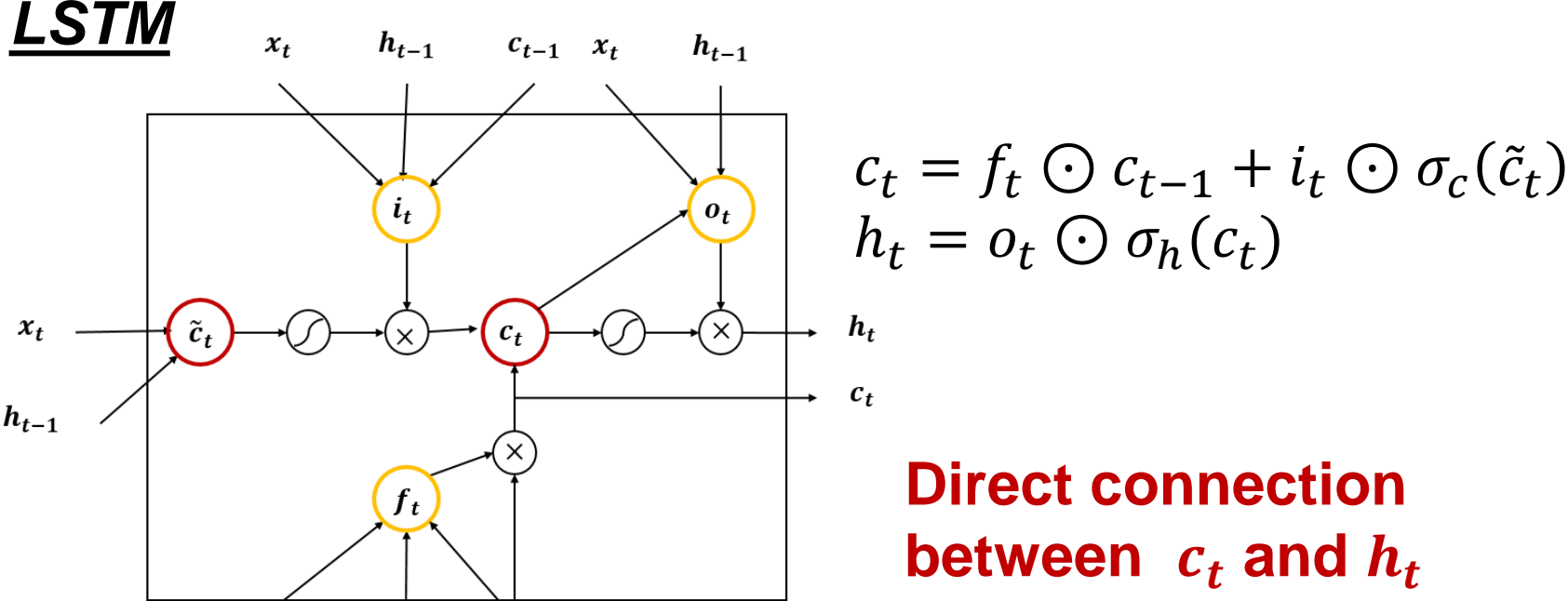
- The bi-channel attentions distinguishes the attentions
- $\alpha_t^{(c)}$  focuses on the neighbor attention (short-term)
- $\alpha_t^{(h)}$  reads out through the whole sequence (long-term)

### 5) Case Study



- $G_{t=17}^{(d)}$  has a relatively small value
- This small value is caused by the selection of different category items to the previous sub-sequence at  $t=16$ .

## Related Works



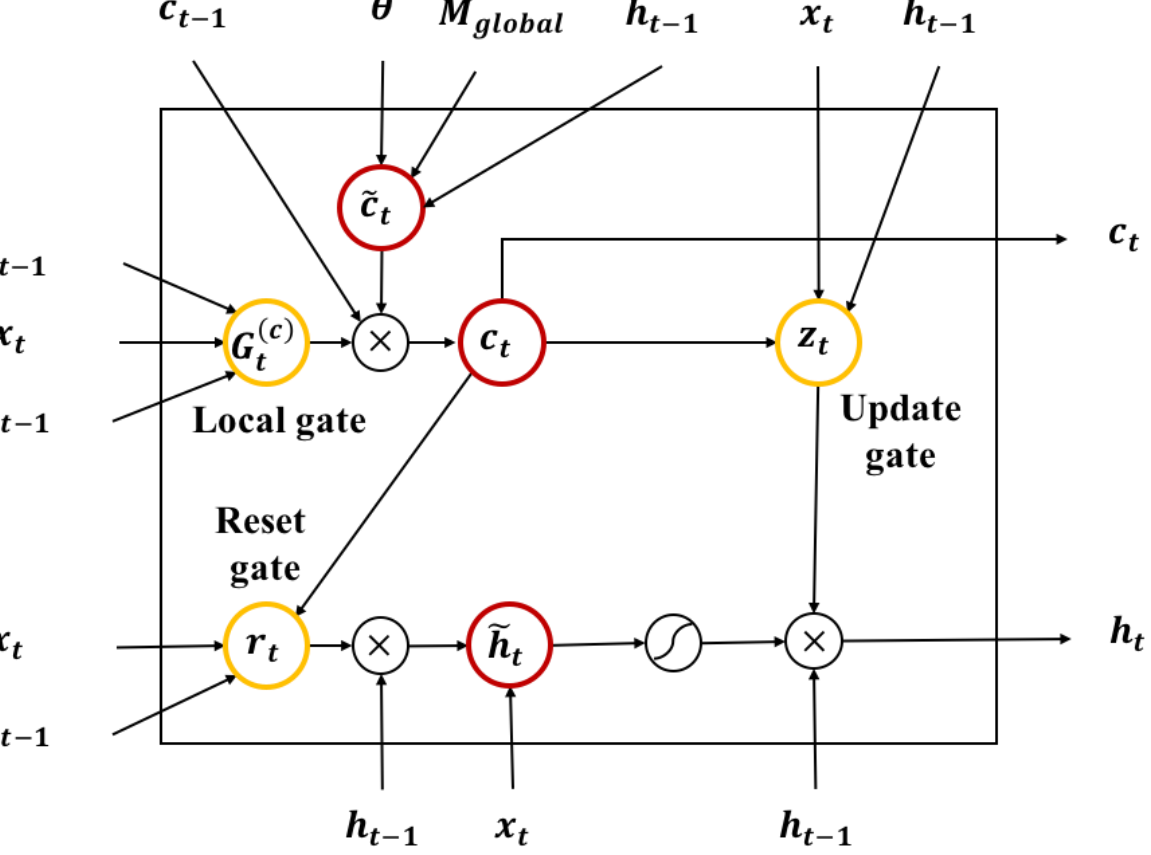
### Sequential Recommendation

- NARM [CIKM-17] : Focus on long-term interest
- STAMP [KDD-18] : Focus on short-term interest
- HCRNN (ours)** : Focus on interest drift with long-term and short-term interest modeling

**There are no studies which capture user's interest change with hierarchical context modeling**

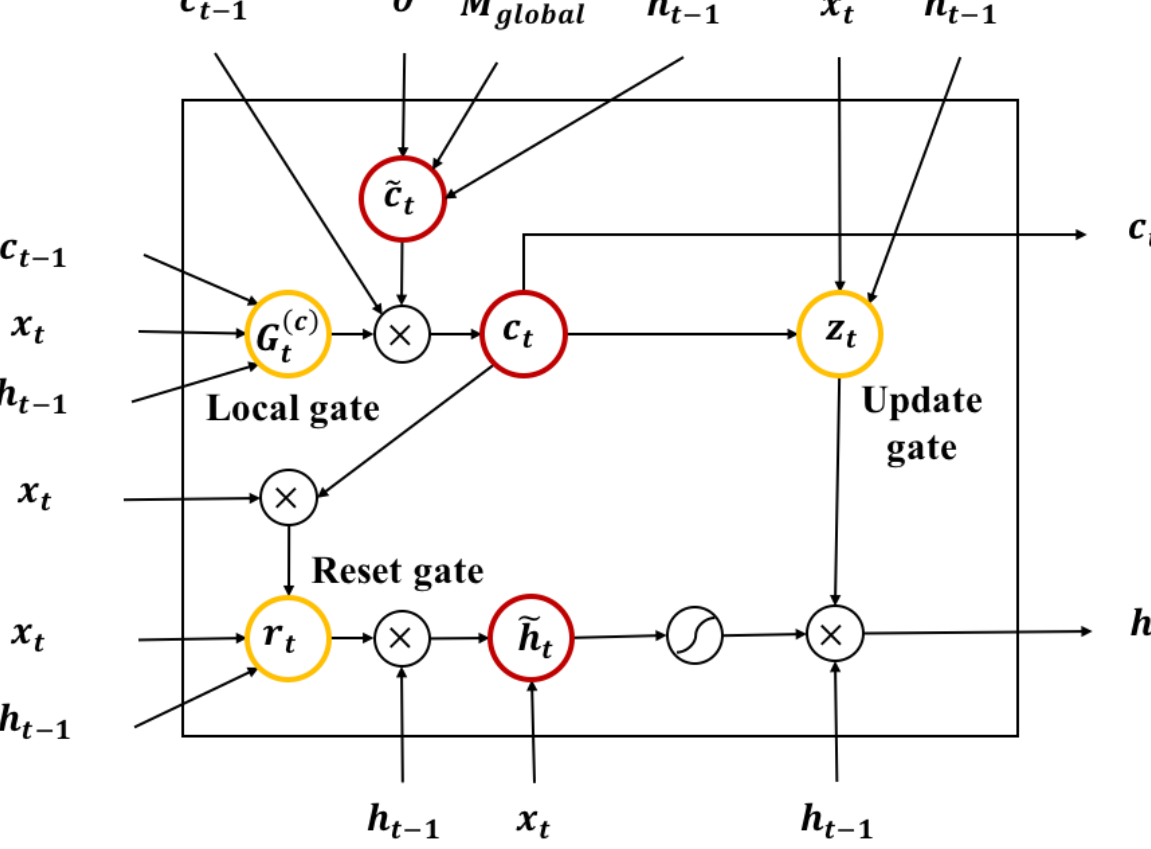
## Methodology

### HCRNN-1



$$\begin{aligned} \tilde{\theta} &\sim q(\tilde{\theta}) = \mathcal{N}(\tilde{\theta}; \mu(x_{1:T}), \text{diag}(\sigma^2(x_{1:T}))) & (11) \\ \theta &\sim \text{softmax}(\tilde{\theta}) & (12) \\ \alpha_t^{(k)} &= \text{softmax}(v_\theta^T \sigma(h_{t-1} W_{h\alpha} + (\theta^{(k)} M_{global}^{(k)}) W_{\theta\alpha})) & (13) \\ \tilde{c}_t &= \sum_{k=1}^K \alpha_t^{(k)} M_{global}^{(k)} & (14) \\ G_t^{(c)} &= \sigma_l(x_t W_{xl} + h_{t-1} W_{hl} + c_{t-1} W_{cl} + b_l) & (15) \\ c_t &= (1 - G_t^{(c)}) \odot c_{t-1} + G_t^{(c)} \odot \tilde{c}_t & (16) \\ z_t &= \sigma_z(x_t W_{xz} + h_{t-1} W_{hz} + c_t W_{zx} + b_z) & (17) \\ r_t &= \sigma_r(x_t W_{xr} + h_{t-1} W_{hr} + c_t W_{cr} + b_r) & (18) \\ \tilde{h}_t &= (r_t \odot h_{t-1}) W_{hh} + x_t W_{xh} + b_h & (19) \\ h_t &= (1 - z_t) \odot h_{t-1} + z_t \odot \sigma_h(\tilde{h}_t) & (20) \end{aligned}$$

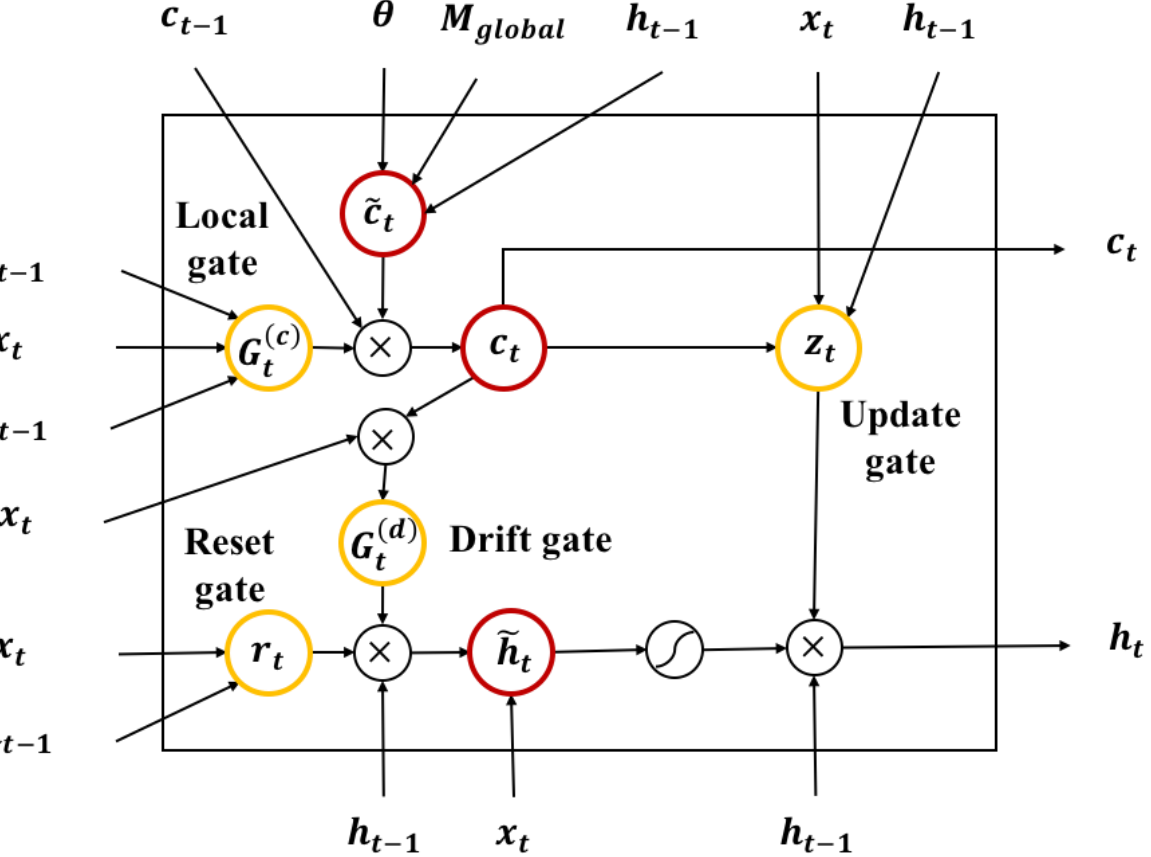
### HCRNN-2



$$\begin{aligned} r_t &= \sigma_r(x_t W_{xr} + h_{t-1} W_{hr} + c_t W_{cr} + b_r) & (18) \\ \tilde{h}_t &= (r_t \odot h_{t-1}) W_{hh} + x_t W_{xh} + b_h & (19) \\ h_t &= (1 - z_t) \odot h_{t-1} + z_t \odot \sigma_h(\tilde{h}_t) & (20) \end{aligned}$$

**Interest drift assumption**  
 $x_t \odot c_t \downarrow \Rightarrow r_t \downarrow \Rightarrow h_t$  focus on the current input instead of  $h_{t-1}$

### HCRNN-3



$$\begin{aligned} r_t &= \sigma_r(x_t W_{xr} + h_{t-1} W_{hr} + c_t W_{cr} + b_r) & (18) \\ \tilde{h}_t &= (r_t \odot h_{t-1}) W_{hh} + x_t W_{xh} + b_h & (19) \\ h_t &= (1 - z_t) \odot h_{t-1} + z_t \odot \sigma_h(\tilde{h}_t) & (20) \end{aligned}$$

$$r_t = \sigma_r(x_t W_{xr} + h_{t-1} W_{hr} + (x_t \odot c_t) W_d + b_r) \quad s.t. W_d \geq 0 \quad (21)$$

Sigmoid function is not sharp  
 $\Rightarrow r_t$  in Eq. 21 :  $0.47 (\pm 0.03)$   
It is hard to incorporate the interest drift

$$\begin{aligned} G_t^{(d)} &= \sigma_d((x_t \odot c_t) W_d + b_d) \quad s.t. W_d \geq 0 & (22) \\ r_t &= \sigma_r(x_t W_{xr} + h_{t-1} W_{hr} + b_r) & (23) \\ \tilde{h}_t &= (r_t \odot (G_t^{(d)} \odot h_{t-1})) W_{hh} + x_t W_{xh} + b_h & (24) \end{aligned}$$

Introduce the interest drift gate ( $G_t^{(d)}$ ) to make  $h_t$  focus on the current input

### HCRNN-3+Bi

$$\begin{aligned} \alpha_{tj}^{(c)} &= \text{softmax}\left(\frac{(c_t W_{c\alpha}^{(1)})(c_j W_{c\alpha}^{(2)})^T}{\sqrt{|H|}}\right) & (25) \\ \alpha_{tj}^{(h)} &= \text{softmax}(v_h^T \sigma(h_t W_{h\alpha}^{(1)} + h_j W_{h\alpha}^{(2)})) & (26) \\ h_t^{(c)} &= \sum_j \alpha_{tj}^{(c)} h_j \quad \text{and} \quad h_t^{(h)} = \sum_j \alpha_{tj}^{(h)} h_j & (27) \\ \hat{y}_t &= \text{softmax}(W_{emb}^T W_B [h_t, h_t^{(c)}, h_t^{(h)}]) & (28) \end{aligned}$$

### Inference

Variational inference by optimizing the evidence lower bound (ELBO)

$$\begin{aligned} \log p(y_{1:T} | c_{1:T}, h_{1:T}) &= \log \int p(\tilde{\theta}) \prod_{t=1}^T p(y_t | \tilde{\theta}, c_t, h_t) d\tilde{\theta} \\ &\geq \sum_{t=1}^T E_{q(\tilde{\theta})} [\log p(y_t | \tilde{\theta}, c_t, h_t)] - \text{KL}([q(\tilde{\theta})] || p(\tilde{\theta})) \end{aligned}$$

- $\alpha_t^{(c)}$  : attention based on the local context (Short-term dependency)
- $\alpha_t^{(h)}$  : attention based on the temporary context (Long-term dependency)