



Pontificia Universidad Católica de Chile  
Escuela de Ingeniería  
Departamento de Ciencia de la Computación  
IIC3633 - Sistemas Recomendadores

# Personalized Transformer for Explainable Recommendation

**Grupo 12:** Daniel Alegría, Gabriel Catalán y Benjamín Faúndez

19/06/2025

**Lei Li<sup>1</sup> Yongfeng Zhang<sup>2</sup> Li Chen<sup>1</sup>**

<sup>1</sup>Hong Kong Baptist University, Hong Kong, China

<sup>2</sup>Rutgers University, New Brunswick, USA

<sup>1</sup>{csleili, lichen}@comp.hkbu.edu.hk

<sup>2</sup>yongfeng.zhang@rutgers.edu

# Resumen

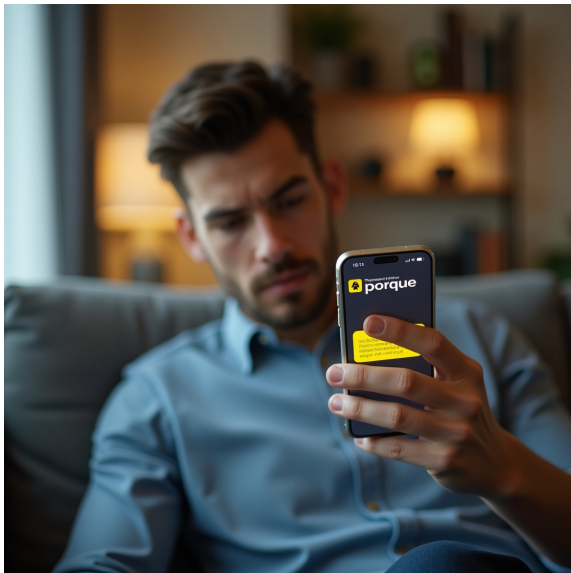
- Contexto
- Problema de recomendación
- Contribución
- Estado del arte y marco teórico
- Detalle solución
- Evaluación
- Referencias
- Preguntas

# Contexto

- v1 Martes, 25 de Mayo 2021
- v2 Sábado, 5 de Junio 2021
- ¿Qué estaba pasando en la escena de la IA?

- Generación de lenguaje natural y sistemas recomendadores
- Usuarios no quieren solo una recomendación, sino que una explicación que justifique la recomendación

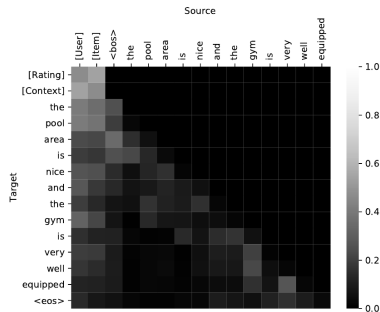
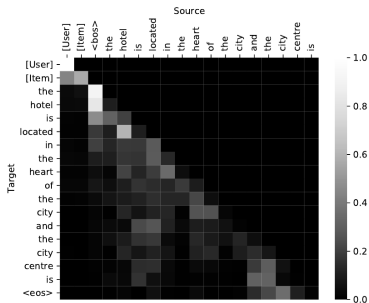
# Contexto: Imaginen...



# Problema de recomendación

- Oración  $\hat{E}_{u,i}$  para un par  $(u, i)$
- PETER también puede estimar  $\hat{r}_{u,i}$
- Se pueden incorporar características de los artículos  $F_{u,i}$

# Limitaciones de los Transformers estándar





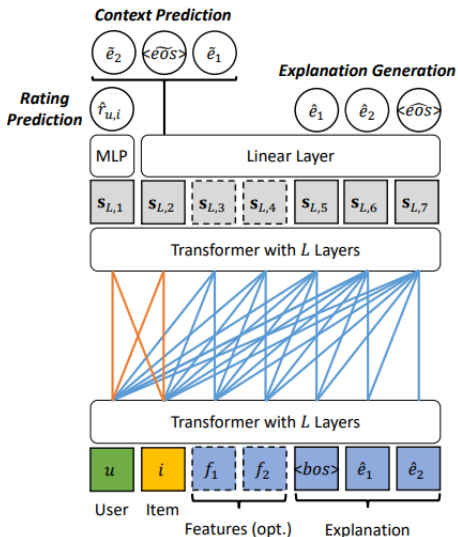
# Contribución

- Según los autores: primer modelo basado en Transformer con generación de lenguaje natural personalizada
- Recomendación y explicación
- Tarea: predicción de contexto
- Pequeño y eficiente

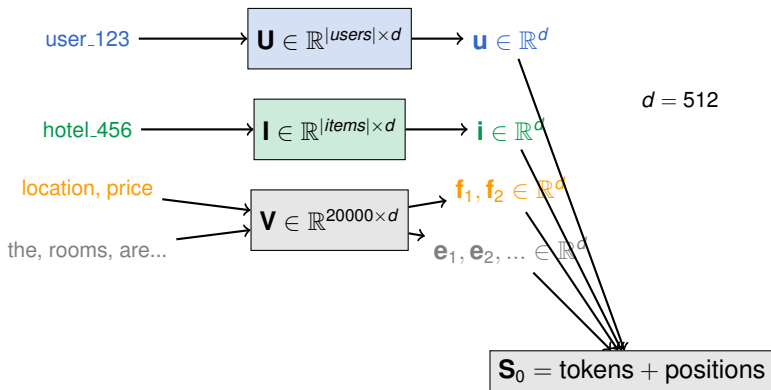
# Estado del arte y marco teórico

- Recomendación explicable: LSTM; Hochreiter and Schmidhuber, 1997. GRU; Cho et al., 2014
- Transformer: Vaswani et al., 2017. Devlin et al., 2019
- Generación personalizada: atributos personales; Zheng et al., 2020. Títulos de películas; Zhou et al., 2020. Características del artículo; Ni et al., 2019

# Solución: Arquitectura PETER



# Representación de Entradas

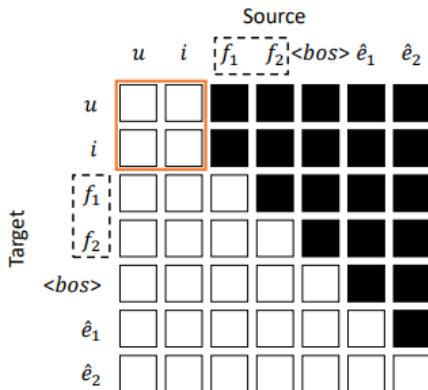


## Ejemplo de secuencia:

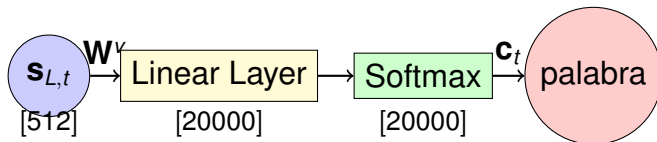
$$S = [u, i, f_1, f_2, \langle \text{bos} \rangle, e_1, e_2, \dots, e_{15}]$$

# Máscara de atención de PETER

 Allow to attend  Prevent from attending



# Generación de Explicación



$$\mathcal{L}_e = -\frac{1}{|T|} \sum_{(u,i) \in T} \frac{1}{|E_{u,i}|} \sum_{t=1}^{|E_{u,i}|} \log c_{2+|F_{u,i}|+t}^{e_t}$$

**Generación paso a paso:**

$\langle \text{bos} \rangle \rightarrow \text{"the"} \rightarrow \text{"hotel"} \rightarrow \text{"is"} \rightarrow \dots \rightarrow \langle \text{eos} \rangle$

# Predicción de Contexto

**Sin Predicción de**

Item ID



“Great location”

Genérico

**Contexto**

**Con Predicción de**

Item ID

pool, nice, gym

“Pool is nice”

Personalizado

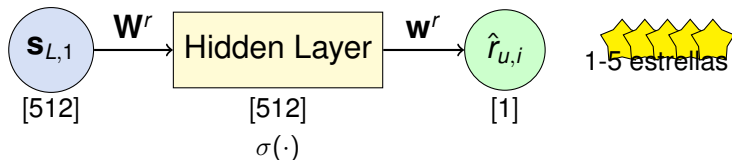
**Contexto**

**Maapeo ID → Palabras**

Posición 2 (ítem) predice TODAS las palabras simultáneamente:

$$\mathcal{L}_c = -\frac{1}{|T|} \sum_{(u,i) \in T} \frac{1}{|E_{u,i}|} \sum_{t=1}^{|E_{u,i}|} \log c_2^{e_t}$$

# Predicción de Rating



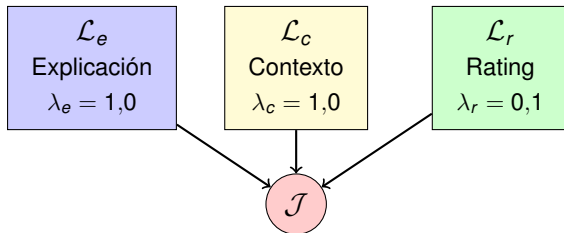
$$\hat{r}_{u,i} = \mathbf{w}^r \sigma(\mathbf{W}^r \mathbf{s}_{L,1} + \mathbf{b}^r) + \bar{b}^r$$

$$\mathcal{L}_r = \frac{1}{|T|} \sum_{(u,i) \in T} (r_{u,i} - \hat{r}_{u,i})^2$$



# Aprendizaje Multitarea

$$\mathcal{J} = \min_{\Theta} (\lambda_e \mathcal{L}_e + \lambda_c \mathcal{L}_c + \lambda_r \mathcal{L}_r)$$



# Evaluación

- Dataset
- Métricas
- Modelos
- Resultados

# Evaluación: Dataset

Se utilizaron tres Datasets

- Yelp
- Amazon
- TripAdvisor

	Yelp	Amazon	TripAdvisor
#users	27,147	7,506	9,765
#items	20,266	7,360	6,280
#records	1,293,247	441,783	320,023
#features	7,340	5,399	5,069
#records / user	47.64	58.86	32.77
#records / item	63.81	60.02	50.96
#words / exp	12.32	14.14	13.01

# Evaluación: Métricas

- RMSE
- MAE
- BLEU-1
- BLEU-4
- ROUGE-1
- ROUGE-2
- USR
- FMR
- FCR
- DIV

# Métricas: BLEU-n

Calcula qué tan parecida es la explicación del modelo con la explicación de testing utilizando n-gram

$$\text{BLEU Score} = BP \cdot \exp \left( \sum_{i=1}^n w_i \cdot \log(p_i) \right) \quad (1)$$

Donde:

- $BP$  corresponde a una penalización si es que la explicación generada por el modelo es más corta que la baseline
- $w_i$  corresponde al peso que se le da a cada i-gram
- $p_i$  es la coincidencia de los i-gram del modelo con la i-gram baseline

# Métricas: Rouge-n

Calcula distintas métricas en base a la superposición de los n-grams de la explicación generada y la explicación testing.

$$\text{Rouge score} = \{ \textit{Precision@n-gram}, \textit{Recall@n-gram}, \textit{F1@n-gram} \} \quad (2)$$

Donde:

- *Precision@n-gram* corresponde a la división entre el numero de n-grams que coinciden en la superposición y el numero de n-grams totales generado por el modelo

# Métricas: Rouge-n

- *Recall@n*-gram corresponde a la división entre el número de n-grams que coinciden en la superposición y el número de n-grams totales del baseline.
- *F1@n*-gram es la métrica usual vista en clases usando las dos métricas anteriores.

# Métricas: Unique Sentence Ratio (USR)

Calcula la división entre el número de explicaciones únicas y el número de explicaciones totales.

$$\text{USR} = \frac{|S|}{N} \quad (3)$$

Donde:

- $S$  es el set de explicaciones únicas
- $N$  es el numero de explicaciones totales.



# Métricas: Feature Matching Ratio (FMR)

Mide cuándo una explicación generada tiene alguna característica de la explicación baseline

$$\text{FMR} = \frac{1}{N} \sum_{u,i} \delta(f_{u,i} \in \hat{S}_{u,i}) \quad (4)$$

Donde:

- $N$  es el numero de explicaciones totales.
- $\hat{S}_{u,i}$  es la explicación generada para algún par usuario-item
- $f_{u,i}$  es la feature dada para ese usuario-item
- $\delta(x) = 1$  si  $x$  es TRUE y  $\delta(x) = 0$  para cualquier otro caso.

# Métricas: Feature Coverage Ratio (FCR)

Mide cuántas características diferentes hay dentro de todas las explicaciones.

$$\text{FCR} = \frac{N_g}{|\mathcal{F}|} \quad (5)$$

Donde:

- $N_g$  es el numero de características diferentes dentro de las explicaciones generadas por el modelo
- $\mathcal{F}$  es el conjunto de todas las características.

# Métricas: Feature Diversity (DIV)

Mide la intersección de las características entre cualquier par de explicaciones generadas.

$$\text{DIV} = \frac{2}{N \cdot (N - 1)} \sum_{u, u', i, i'} \left| \hat{\mathcal{F}}_{u,i} \cap \hat{\mathcal{F}}_{u',i'} \right| \quad (6)$$

Donde:

- $N$  es el numero de explicaciones totales.
- $\hat{\mathcal{F}}_{u,i}$  son el conjunto de las características presentes en una explicación generada por el modelo para el usuario  $u$  y para el item  $i$
- $\hat{\mathcal{F}}_{u',i'}$  son el conjunto de las características presentes en una explicación generada por el modelo para el usuario  $u'$  y para el item  $i'$

# Evaluación: Modelos

## Modelos de explicación

- Transformer
- NRT
- Att2Seq
- PETER
- \*ACMLM
- \*NETE
- \*PETER+

## Modelos de recomendación

- NRT
- PETER
- \*NETE
- PMF
- SVD++

Todos los modelos con \* aceptan características como input para el entrenamiento.

# Evaluación: Resultados de explicabilidad

	Explainability			Text Quality								
	FMR↑	FCR↑	DIV↓	USR↑	B1↑	B4↑	R1-P↑	R1-R↑	R1-F↑	R2-P↑	R2-R↑	R2-F↑
Yelp												
Transformer	0.06	0.06	2.46	0.01	7.39	0.42	<b>19.18</b>	10.29	12.56	1.71	0.92	1.09
NRT	0.07	0.11	2.37	0.12	<b>11.66</b>	0.65	17.69	12.11	13.55	1.76	1.22	1.33
Att2Seq	0.07	<u>0.12</u>	2.41	<b>0.13</b>	10.29	0.58	<u>18.73</u>	11.28	13.29	<u>1.85</u>	1.14	1.31
PETER	<b>0.08**</b>	<b>0.19**</b>	<b>1.54**</b>	<b>0.13</b>	10.77	<b>0.73**</b>	18.54	<b>12.20</b>	<b>13.77**</b>	<b>2.02**</b>	<b>1.38**</b>	<b>1.49**</b>
ACMLM	0.05	<u>0.31</u>	<b>0.95</b>	<b>0.95</b>	7.01	0.24	7.89	7.54	6.82	0.44	0.48	0.39
NETE	0.80	0.27	1.48	0.52	<u>19.31</u>	<u>2.69</u>	<u>33.98</u>	<u>22.51</u>	25.56	8.93	5.54	6.33
PETER+	<b>0.86**</b>	<b>0.38**</b>	<u>1.08</u>	0.34	<b>20.80**</b>	<b>3.43**</b>	<b>35.44**</b>	<b>26.12**</b>	<b>27.95**</b>	<b>10.65**</b>	<b>7.44**</b>	<b>7.94**</b>
Amazon												
Transformer	0.10	0.01	3.26	0.00	9.71	0.59	19.68	11.94	14.11	2.10	1.39	1.55
NRT	<b>0.12</b>	0.07	2.93	0.17	<b>12.93</b>	0.96	<b>21.03</b>	13.57	<b>15.56</b>	<u>2.71</u>	1.84	2.05
Att2Seq	<b>0.12</b>	<u>0.20</u>	<u>2.74</u>	<b>0.33</b>	12.56	0.95	20.79	13.31	<u>15.35</u>	2.62	1.78	1.99
PETER	<b>0.12**</b>	<b>0.21</b>	<b>1.75**</b>	0.29	<u>12.77</u>	<b>1.17**</b>	19.81	<b>13.80</b>	15.23	<b>2.80</b>	<b>2.08**</b>	<b>2.20**</b>
ACMLM	0.10	<b>0.31</b>	2.07	<b>0.96</b>	9.52	0.22	11.65	10.39	9.69	0.71	0.81	0.64
NETE	<u>0.71</u>	0.19	<u>1.93</u>	0.57	18.76	2.46	<u>33.87</u>	<u>21.43</u>	<u>24.81</u>	<u>7.58</u>	4.77	5.46
PETER+	<b>0.77**</b>	<b>0.31**</b>	<b>1.20**</b>	0.46	<b>19.75**</b>	<b>3.06**</b>	<b>34.71**</b>	<b>23.99**</b>	<b>26.35**</b>	<b>9.04**</b>	<b>6.23**</b>	<b>6.71**</b>
TripAdvisor												
Transformer	0.04	0.00	10.00	0.00	12.79	0.71	16.52	<b>16.38</b>	15.88	2.22	<b>2.63</b>	<b>2.34</b>
NRT	<u>0.06</u>	0.09	<u>4.27</u>	0.08	15.05	0.99	18.22	14.39	15.40	2.29	1.98	2.01
Att2Seq	<u>0.06</u>	<b>0.15</b>	4.32	<b>0.17</b>	<u>15.27</u>	1.03	18.97	14.72	<u>15.92</u>	<b>2.40</b>	2.03	2.09
PETER	<b>0.07**</b>	0.13	<b>2.95**</b>	0.08	<b>15.96**</b>	<b>1.11*</b>	<b>19.07</b>	<u>16.09</u>	<b>16.48**</b>	<u>2.33</u>	2.17	<u>2.09</u>
ACMLM	0.07	<b>0.41</b>	<b>0.78</b>	<b>0.94</b>	3.45	0.02	4.86	3.82	3.72	0.18	0.20	0.16
NETE	<u>0.78</u>	0.27	2.22	0.57	<u>22.39</u>	<u>3.66</u>	<u>35.68</u>	<u>24.86</u>	<u>27.71</u>	<u>10.20</u>	6.98	<u>7.66</u>
PETER+	<b>0.89**</b>	<u>0.35</u>	<u>1.61</u>	0.25	<b>24.32**</b>	<b>4.55**</b>	<b>37.48**</b>	<b>29.21**</b>	<b>30.49**</b>	<b>11.92**</b>	<b>8.98**</b>	<b>9.24**</b>

# Evaluación: Análisis de tiempo

	Time	Epochs	Time/Epoch
ACMLM	97.0	<b>3</b>	32.3
PETER+	<b>57.7</b>	25	<b>2.3</b>

# Evaluación: Análisis cualitativo

	Top-15 Context Words	Explanation
Ground-truth		
PETER	<eos> the and a <u>pool</u> was with nice is very were to good in of	the <b>rooms</b> are spacious and the bathroom has a large tub
PETER+	<eos> the and a was <u>pool</u> with to nice good very were is of in	the <u>pool</u> area is nice and the <u>gym</u> is very well equipped <eos>
Ground-truth		
PETER	<eos> the and a was were separate bathroom with <u>shower</u> large very had in is	the <u>rooms</u> were clean and comfortable <eos>
PETER+	<eos> the and a was <u>bathroom</u> <u>shower</u> with large in separate were <u>room</u> very is	beautiful <b>lobby</b> and nice bar
		the bathroom was large and the <u>shower</u> was great <eos>
		the <u>lobby</u> was very nice and the <u>rooms</u> were very comfortable <eos>

# Evaluación: Resultados de recomendación

	Yelp		Amazon		TripAdvisor	
	R↓	M↓	R↓	M↓	R↓	M↓
PMF	1.09	0.88	1.03	0.81	0.87	0.70
SVD++	<b>1.01</b>	<b>0.78</b>	0.96	0.72	0.80	0.61
NRT	<b>1.01</b>	<b>0.78</b>	<b>0.95</b>	<b>0.70</b>	<b>0.79</b>	0.61
NETE	<b>1.01</b>	0.79	0.96	0.73	<b>0.79</b>	<b>0.60</b>
PETER	<b>1.01</b>	<b>0.78</b>	<b>0.95</b>	0.71	0.81	0.63



# Evaluación: Análisis de ablación

	Explainability			Text Quality			Recommendation	
	FMR	FCR	DIV	USR	BLEU-1	BLEU-4	RMSE	MAE
Disable $\mathcal{L}_c$	0.06 ↓	0.03 ↓	5.75 ↓	0.01 ↓	15.37 ↓	0.86 ↓	0.80 ↑	0.61 ↑
Disable $\mathcal{L}_r$	0.07	0.14 ↑	2.90 ↑	0.10 ↑	16.16 ↑	1.15 ↑	3.23 ↓	3.10 ↓
Left-to-Right Masking	0.07	0.15 ↑	2.68 ↑	0.12 ↑	15.73 ↓	1.11	0.87 ↓	0.68 ↓
PETER	0.07	0.13	2.95	0.08	15.96	1.11	0.81	0.63

# Referencias

- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding.

# Referencias

- Yinhe Zheng, Rongsheng Zhang, Minlie Huang, and Xiaoxi Mao. 2020. A pre-training based personalized dialogue generation model with persona-sparse data.
- Kun Zhou, Wayne Xin Zhao, Shuqing Bian, Yuanhang Zhou, Ji-Rong Wen, and Jingsong Yu. 2020. Improving conversational recommender systems via knowledge graph based semantic fusion.
- Piji Li, Zihao Wang, Zhaochun Ren, Lidong Bing, and Wai Lam. 2017. Neural rating regression with abstractive tips generation for recommendation.
- Li Dong, Shaohan Huang, Furu Wei, Mirella Lapata, Ming Zhou, and Ke Xu. 2017. Learning to generate product reviews from attributes.

# Referencias

- Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects.
- Lei Li, Yongfeng Zhang, and Li Chen. 2020c. Generate neural template explanations for recommendation.
- Andriy Mnih and Russ R Salakhutdinov. 2007. Probabilistic matrix factorization.
- Yehuda Koren. 2008. Factorization meets the neighborhood: a multifaceted collaborative filtering model.

# Preguntas

