

Rethinking Positional Encoding in Tree Transformer for Code Representation



Han Peng, Ge Li*, Yunfei Zhao, Zhi Jin*

Key Laboratory of High Confidence Software Technologies (Peking University), Ministry of Education;
Institute of Software, EECS, Peking University, Beijing, China

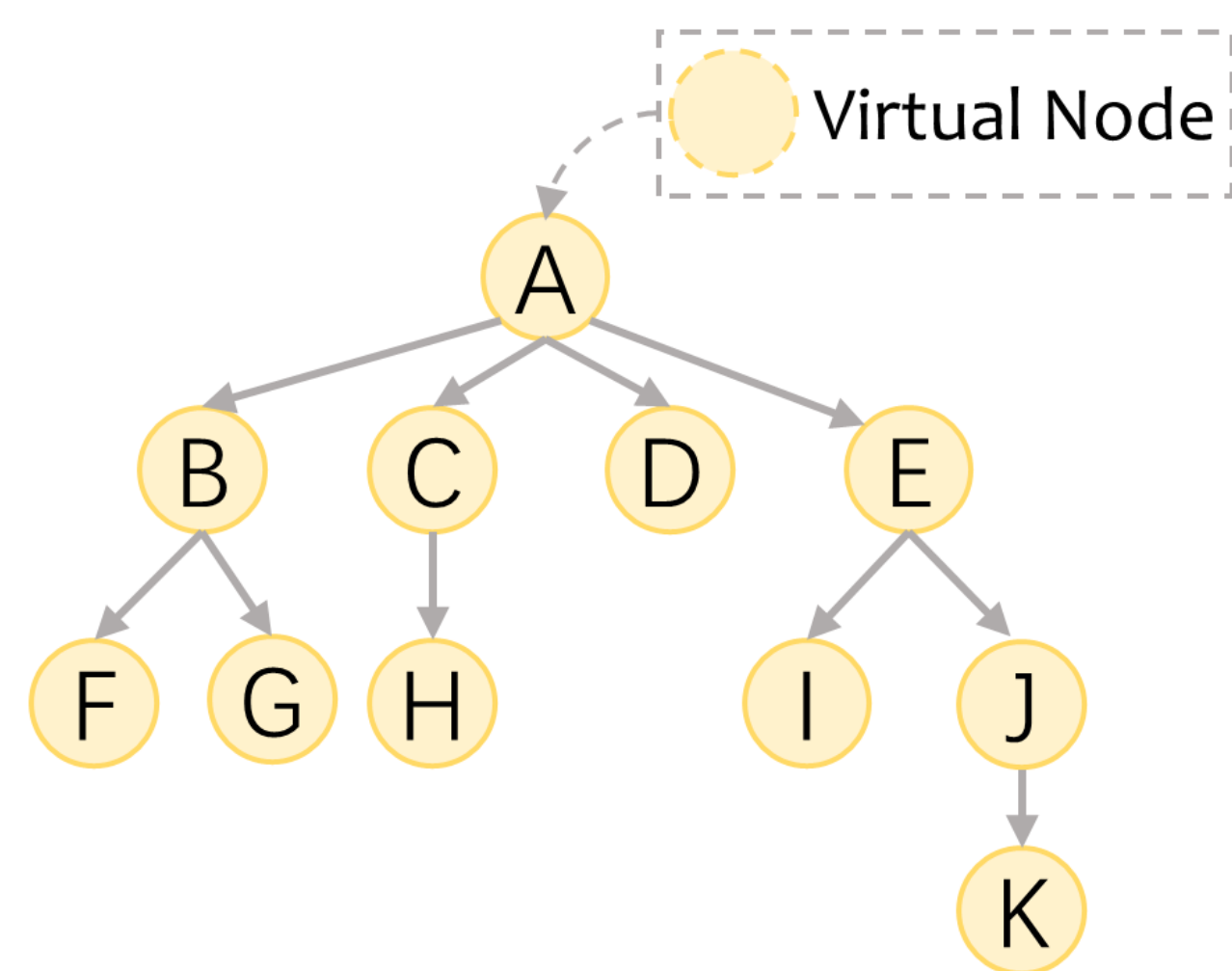
Introduction

Several recent works develop tree Transformers to capture the AST structure in source code, and specifically, novel tree positional encodings have been proposed to incorporate structure inductive bias into Transformer.

In this work, we propose a novel tree Transformer encoding node positions based on our new description method for tree structures. Technically, local and global soft bias is both introduced as positional encodings of our model.

Approach

1) A 2D recursive description for code AST



- A : $\{(1,1)\}$
- B : $\{(1,1), (1,4)\}$
-
- J : $\{(1,1), (4,4), (2,2)\}$
- K : $\{(1,1), (4,4), (2,2), (1,1)\}$

Our proposed description for tree structure is defined from the tree root to leaves recursively. Specifically, the position for each node is represented as:

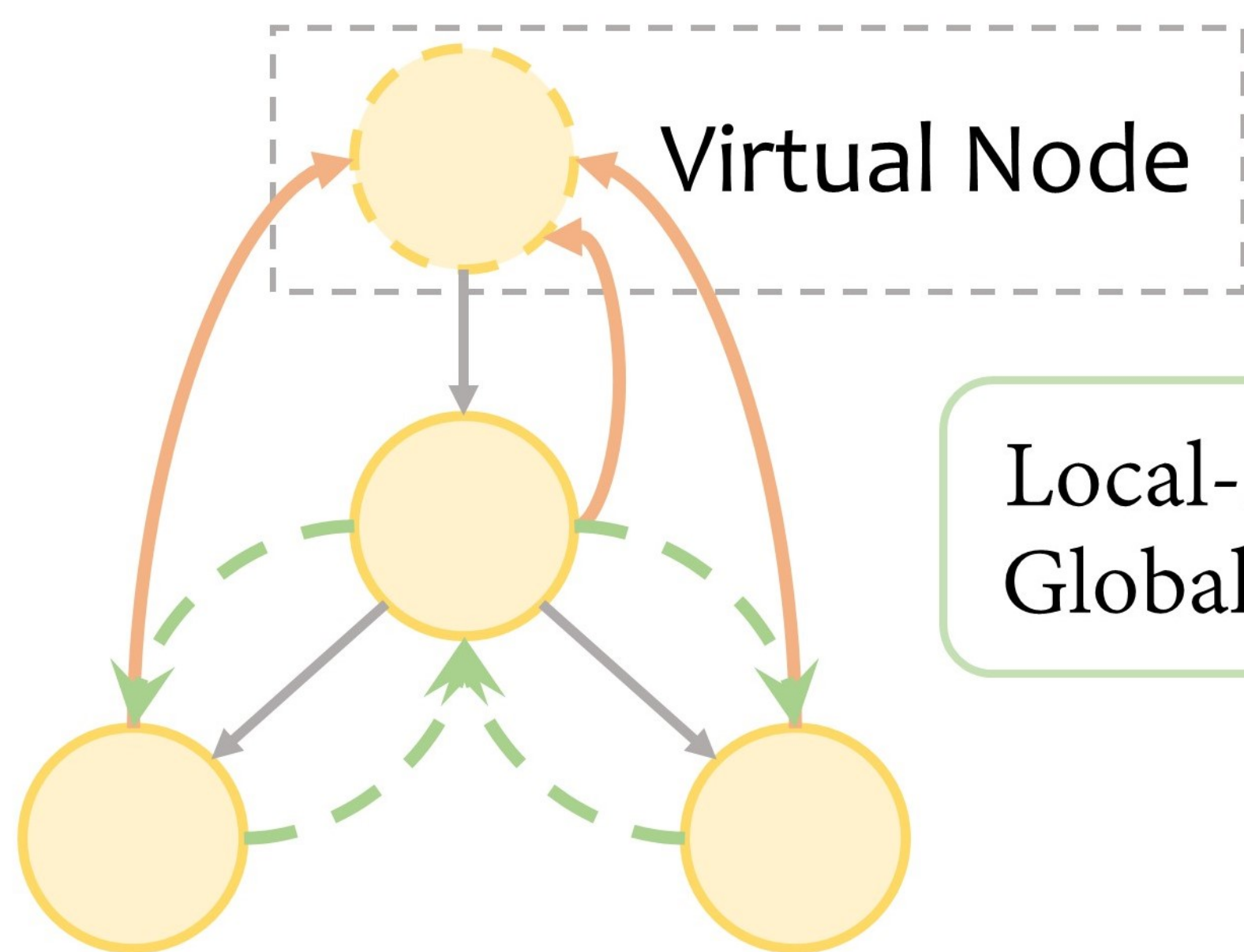
$$\mathcal{F}(x) = \begin{cases} \mathcal{F}(f(x)) + \{(x_i, x_j)\} & \text{if } x \neq \text{root} \\ \{(1,1)\} & \text{if } x = \text{root} \end{cases}$$

Emb

$$H_i = [h(i^1), h(i^2), \dots, h(i^n)]$$

The $\mathcal{F}(x)$ is the position description for node x and $f(x)$ is specified as the parent node for it. The x_i is the sibling order of node x and the x_j is the total child number of its parent $f(x)$. The n is the depth of node i in the tree and $h(i^n)$ is the embedding vector of the n th two-dimensional coordinate.

2) Encoding tree positions in local and global



Local-rel:
 Global-abs:

$$r_{ij} = \begin{cases} LN(Linear(\sum_h H_i - \sum_h H_j)) & \text{if } f(i) = j \vee f(j) = i \\ \vec{0} & \text{if } f(i) \neq j \wedge f(j) \neq i \end{cases}$$

$$\gamma_{ij} = (x_i W^Q)(r_{ij} W_r^K)^T + (r_{ji} W_r^Q)(x_j W^K)^T$$

Local-rel

$$a_i = LN(Linear(Concat(H_i)))$$

$$\beta_{ij} = (a_i W_a^Q)(a_j W_a^K)^T$$

Global-abs

$$A_{ij} = \frac{1}{\sqrt{2}}(\alpha_{ij} + \beta_{ij} + \gamma_{ij}),$$

$$z_i = \sum_{j=1}^n \frac{\exp(A_{ij})}{\sum_{j'=1}^n \exp(A_{ij'})} (x_j W^V)$$

$$\alpha_{ij} = (x_i W^Q)(x_j W^K)^T$$

We feed all nodes into Transformer and integrate positional vectors into self-attention as structure bias by positional encoding.

Experiment

Model	Python					JavaScript				
	MRR (Type)	MRR (Value)	ACC (Type)	ACC (Value)	ACC (All)	MRR (Type)	MRR (Value)	ACC (Type)	ACC (Value)	ACC (All)
Transformer	88.76	54.28	81.91	49.56	59.27	89.51	62.73	82.73	57.40	63.34
(Shaw et al., 2018)	± 0.09	± 0.02	± 0.13	± 0.06	± 0.11	± 0.01	± 0.02	± 0.01	± 0.02	± 0.03
(Shiv and Quirk, 2019)	89.05	55.11	82.34	50.48	60.07	90.07	64.13	83.60	59.09	64.80
(Kim et al., 2020)	± 0.04	± 0.04	± 0.06	± 0.08	± 0.14	± 0.03	± 0.05	± 0.05	± 0.09	± 0.14
Our model	91.58	55.62	86.15	51.19	63.19	91.63	64.08	85.91	58.99	66.48
w/o first dim.	± 0.11	± 0.09	± 0.17	± 0.14	± 0.26	± 0.06	± 0.22	± 0.10	± 0.28	± 0.29
w/o second dim.	90.78	55.10	84.87	50.52	62.04	91.08	63.50	85.02	58.30	65.47
w/o Global-Abs.	91.40	55.25	85.88	50.68	62.70	91.57	63.95	85.84	58.84	66.34
w/o Local-Rel.	± 0.07	± 0.03	± 0.12	± 0.08	± 0.15	± 0.07	± 0.15	± 0.11	± 0.18	± 0.12

Table 1: All results on code completion

Model	Python			JavaScript		
	Bleu	F1	Acc	Bleu	F1	Acc
Transformer	56.29 \pm 0.06	33.95 \pm 0.28	18.29 \pm 0.44	56.51 \pm 0.13	22.62 \pm 0.59	11.66 \pm 0.34
(Shaw et al., 2018)	56.86 \pm 0.30	35.58 \pm 0.48	19.51 \pm 0.41	56.75 \pm 0.13	23.22 \pm 0.73	11.32 \pm 1.78
(Shiv and Quirk, 2019)	56.50 \pm 0.19	34.54 \pm 0.38	18.89 \pm 0.48	57.04 \pm 0.23	24.01 \pm 0.21	11.53 \pm 0.53
(Kim et al., 2020)	57.02 \pm 0.17	35.55 \pm 0.41	19.29 \pm 0.35	57.33 \pm 0.38	23.47 \pm 0.28	12.07 \pm 0.74
Our model	57.28\pm0.16	36.00\pm0.10	19.97\pm0.17	57.72\pm0.19	25.03\pm0.37	13.21\pm0.24
w/o first dim.	57.10 \pm 0.14	35.94 \pm 0.21	19.54 \pm 0.20	57.40 \pm 0.25	24.22 \pm 0.65	12.92 \pm 0.27
w/o second dim.	56.89 \pm 0.15	35.36 \pm 0.10	19.24 \pm 0.48	57.27 \pm 0.09	24.30 \pm 0.29	12.66 \pm 0.69
w/o Global-Abs.	56.83 \pm 0.14	35.80 \pm 0.29	19.47 \pm 0.10	57.27 \pm 0.32	23.86 \pm 0.10	12.30 \pm 0.65
w/o Local-Rel.	56.95 \pm 0.10	35.36 \pm 0.37	19.29 \pm 0.66	57.50 \pm 0.05	23.74 \pm 0.18	13.13 \pm 0.09

Table 2: All results on code summarization

We focus on two tasks of code representation learning: code completion and summarization. Our model substantially outperforms baselines almost in all metrics of both tasks. Besides, we find that both two dimensions in the coordinate list are helpful and also confirm the benefits of fusing the local and global paradigms.