

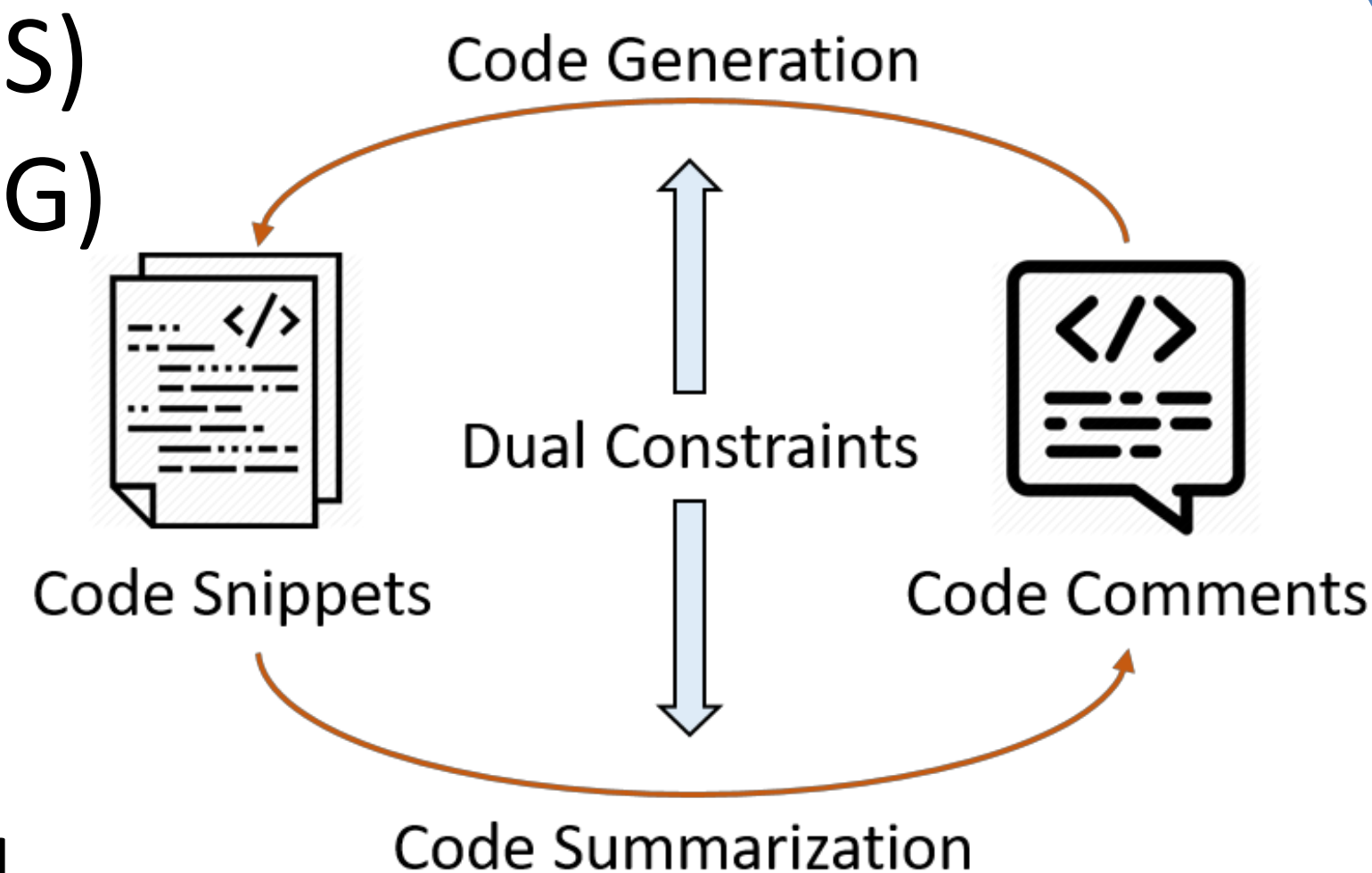
# Code Generation as a Dual Task of Code Summarization

Bolin Wei<sup>1</sup>, Ge Li<sup>1</sup>, Xin Xia<sup>2</sup>, Zhiyi Fu<sup>1</sup>, Zhi Jin<sup>1</sup>

<sup>1</sup>Key Laboratory of High Confidence Software Technologies (Peking University) Ministry of Education, China;  
<sup>1</sup>Software Institute, Peking University, China <sup>2</sup>Faculty of Information Technology, Monash University, Australia

## Introduction

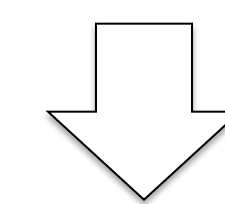
Code summarization (CS) and code generation (CG) are two crucial tasks in the field of software development. None of the previous studies before have considered the relations between the two tasks or exploited the relations to improve each other. Thus, we design a dual learning framework to train a CS and a CG model simultaneously to exploit the duality of them. Besides a probabilistic correlation, we design a novel dual constraint about attention mechanism.



## Dual Constraints

### Probability:

$$p(x|y) * p(y) = p(y|x) * p(x)$$



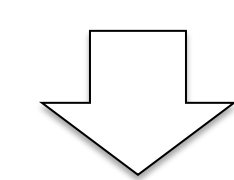
Lagrange multipliers

$$\log p(x) + \log p(y|x) - \log p(x|y) - \log p(y)$$

### Attention Weights:

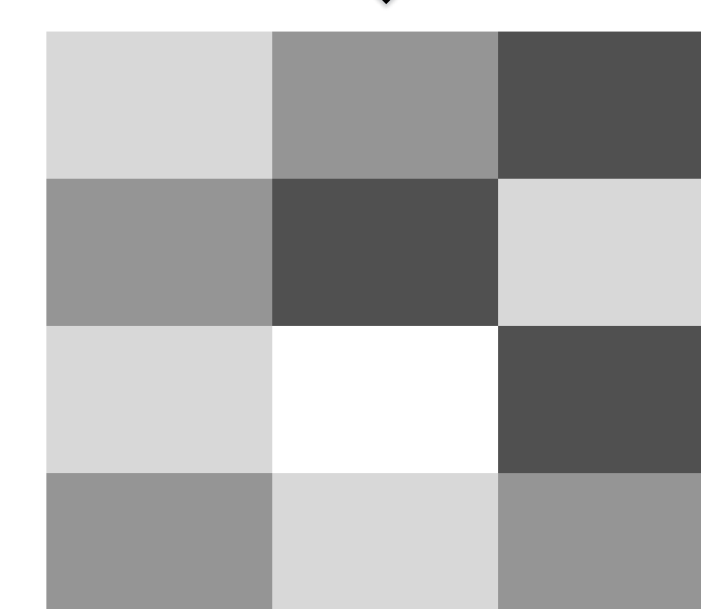
$$b_i = \text{softmax}(A_{xy}^i) \text{ from CS}$$

$$b'_i = \text{softmax}(A_{yx}^i) \text{ from CG}$$



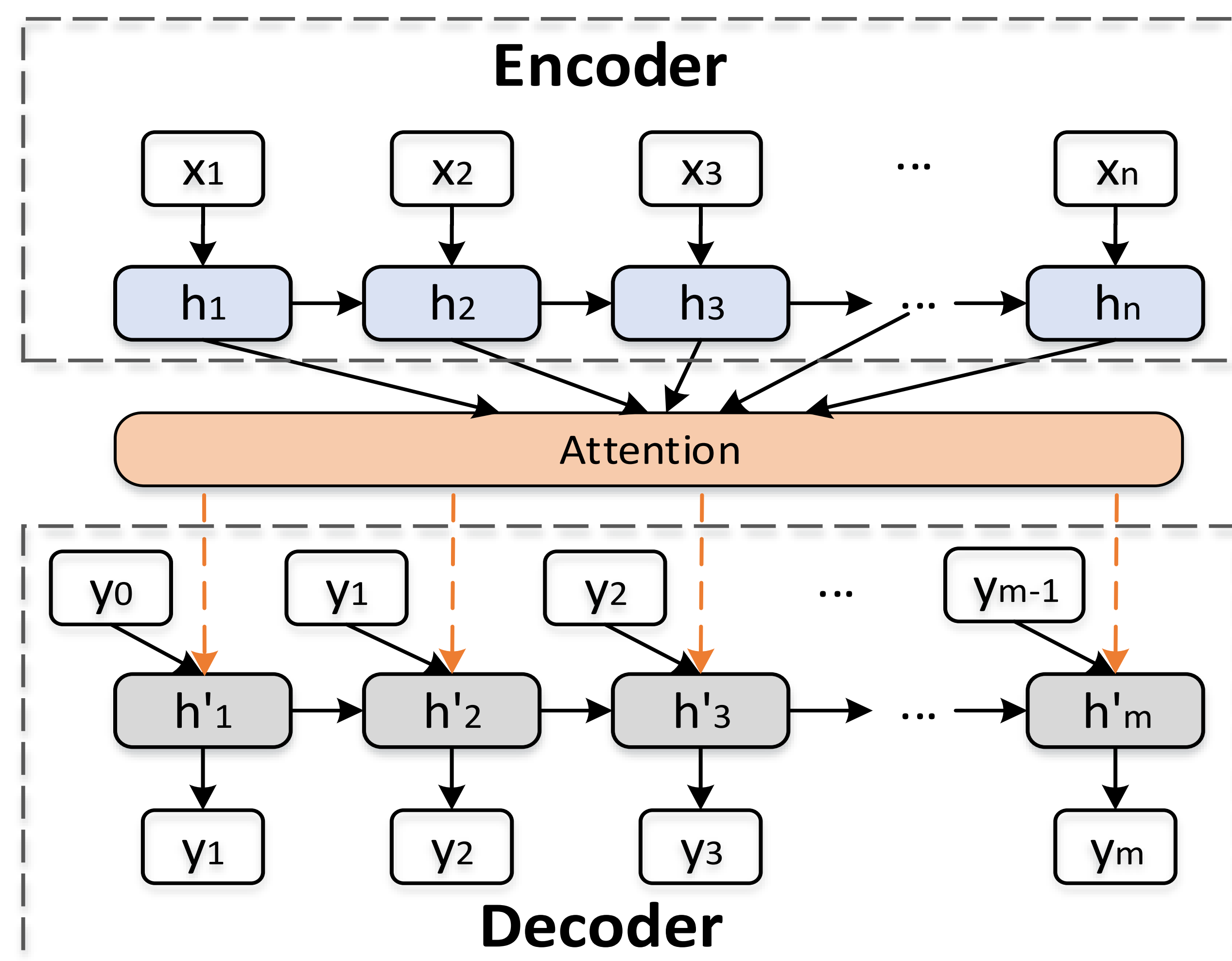
Using KL divergence to constrain the distance

i-th token in comment



j-th token in code

## CS & CG Model



## Algorithm

### Algorithm 1 Algorithm Description

**Input:** Language models  $\hat{P}(x)$  and  $\hat{P}(y)$  for any  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ ;  
hyper parameters  $\lambda_{dual1}$ ,  $\lambda_{dual2}$ ,  $\lambda_{att1}$  and  $\lambda_{att2}$ ;  
optimizers  $opt1$  and  $opt2$

**repeat**

Get a minibatch of  $k$  pairs  $\langle (x_i, y_i) \rangle_{i=1}^k$ ;

Calculate the gradients for  $\theta_{xy}$  and  $\theta_{yx}$ .

$$G_{xy} = \nabla_{\theta_{xy}} (1/k) \sum_{i=1}^k [l_{xy} + \lambda_{dual1} l_{dual} + \lambda_{att1} l_{att}];$$

$$G_{yx} = \nabla_{\theta_{yx}} (1/k) \sum_{i=1}^k [l_{yx} + \lambda_{dual2} l_{dual} + \lambda_{att2} l_{att}];$$

Update  $\theta_{xy}$  and  $\theta_{yx}$

$\theta_{xy} \leftarrow opt1(\theta_{xy}, G_{xy})$ ,  $\theta_{yx} \leftarrow opt2(\theta_{yx}, G_{yx})$

**until** models converged

## CS Results

Table 2: The overall performance of our CS models compared with baselines

Methods	Java			Python		
	BLEU	METEOR	ROUGE-L	BLEU	METEOR	ROUGE-L
CODE-NN	27.60	12.61	41.10	17.36	9.288	37.81
DeepCom	39.75	23.06	52.67	20.78	9.979	37.35
Tree2Seq	37.88	22.55	51.50	20.07	8.957	35.64
RL+Hybrid2Seq	38.22	22.75	51.91	19.28	9.752	39.34
API+CODE	41.31	23.73	52.25	15.36	8.571	33.65
Basic Model	41.01	23.26	51.64	20.47	10.38	38.77
Dual Model	<b>42.39</b>	<b>25.77</b>	<b>53.61</b>	<b>21.80</b>	<b>11.14</b>	<b>39.45</b>

Our dual model obviously outperforms all the baselines on three metrics at the same time. The results of Wilcoxon Rank Sum test and Cliff's delta both indicate the improvements are significant.

## CG Results

Table 3: BLEU scores and percentage of valid code (PoV) on CG task

Methods	Java		Python	
	BLEU	PoV	BLEU	PoV
SEQ2TREE	13.80	22.6%	4.472	22.7%
Basic Model	10.86	19.6%	10.43	41.8%
Dual Model	<b>17.17</b>	<b>27.4%</b>	<b>12.09</b>	<b>51.9%</b>

Although the BLEU score is very low, dual training can still improve the performance of individually trained basic model, proving the effectiveness of dual training.

## Conclusion

We aim to build a framework which uses the CG as a dual task for the CS. We applied two constraints on the dual training loss function. Experimental results show that the dual training process helps the CS and CG model surpass the existing SOTA methods on two datasets.