

supplement of XYLayoutLM: Towards Layout-Aware Multimodal Networks For Visually-Rich Document Understanding

Zhangxuan Gu^{1,2}, Changhua Meng², Ke Wang², Jun Lan², Weiqiang Wang², Ming Gu², Liqing Zhang^{1*}

¹MoE Key Lab of Artificial Intelligence, Shanghai Jiao Tong University, ²Ant Group

zhangxgu@126.com

{changhua.mch, kaywang.wk, yelan.lj, weiqiang.wwq, guming.mg}@antgroup.com

zhang-lq@cs.sjtu.edu.cn

Method	Inference time/one doc	Device
LayoutReader	(10.3ms \pm 34.1us)*1024	one V100
XY Cut	8.99ms \pm 28.3 μ s	CPU

Table 1. Inference time for detecting reading order.

1. Appendix

1.1. Pseudo-code

The pseudo-code is shown in Algorithm 1.

1.2. Inference time for LayoutReader and our method

We show the inference time for our XYLayoutLM and LayoutReader in Table 1. The *1024 means it need to decode 1024 times for one document by default.

*Corresponding author.

Algorithm 1 Augmented XY Cut Algorithm

Require: boxes: $B = \{b_i\}_{i=1}^K$, thresholds: $\lambda_x, \lambda_y, \theta$

Ensure: proper reading order: $O = \{s(i)\}_{i=1}^K$

```
1: function CUT(boxes, n, result, tmp, direction)
2:   if  $\text{len}(\text{boxes}) = 0$  or  $n$  then
3:     return result
4:   end if
5:   sort boxes by direction           ▷ also sort tmp
6:   if direction is Y-axis then
7:     next  $\leftarrow$  X-axis
8:   else if direction is X-axis then
9:     next  $\leftarrow$  Y-axis
10:  end if
11:  cur  $\leftarrow$  0
12:  sets  $\leftarrow$  project boxes to direction
13:  for  $i$  in  $\text{range}(\text{len}(\text{boxes}))$  do
14:    set  $\leftarrow$  sets[ $i$ ]
15:    if  $\text{set} \cap \text{sets}[i:] = \emptyset$  then
16:      result += CUT(boxes[cur :  $i + 1$ ],  $i -$ 
17:      cur, [], tmp[cur :  $i + 1$ ], next)
18:      cur  $\leftarrow$   $i + 1$ 
19:    end if
20:  end for
21:  if  $\text{cur} \neq i + 1$  then
22:    result += tmp[cur :  $i + 1$ ]
23:  end if
24:  return result
25: end function
26: tmp  $\leftarrow$   $\text{range}(K)$ 
27: for  $i$  in  $\text{range}(\text{len}(B))$  do
28:    $b \leftarrow B[i]$ 
29:   Random init  $v_x \in N(-1, 1), v_y \in N(-1, 1)$ 
30:   if  $|v_x| > \lambda_x$  then
31:      $b[0] += \theta \cdot v_x, b[2] += \theta \cdot v_x$ 
32:   end if
33:   if  $|v_y| > \lambda_y$  then
34:      $b[1] += \theta \cdot v_y, b[3] += \theta \cdot v_y$ 
35:   end if
36: end for
37:  $O \leftarrow$  CUT( $B, K, [], \text{tmp}, \text{Y-axis}$ )
```
