

**Algorithm 1** Deterministic Hybrid Layout-OCR Fusion

---

**Require:** Set of AWS layout blocks  $B$ ; ordered list of Google OCR word tokens  $T$ ; tolerance  $\tau = 0.02$

**Ensure:** Structured document object  $D$

```

1: for all  $b \in B$  do
2:    $b.\text{coords} \leftarrow [\text{Left}(b), \text{Top}(b), \text{Right}(b), \text{Bottom}(b)]$ 
3:    $b.\text{area} \leftarrow (\text{Right}(b) - \text{Left}(b)) \times (\text{Bottom}(b) - \text{Top}(b))$ 
4:    $b.\text{matched} \leftarrow \emptyset$ 
5: end for
6: Sort  $B$  in ascending order of  $b.\text{area}$ 
7:  $U \leftarrow \emptyset$  ▷ Used word indices
8: for  $i \leftarrow 1$  to  $|T|$  do
9:    $(x, y) \leftarrow \text{Centroid}(T[i])$ 
10:  for all  $b \in B$  do
11:    if  $\text{Left}(b) - \tau \leq x \leq \text{Right}(b) + \tau$  and  $\text{Top}(b) - \tau \leq y \leq \text{Bottom}(b) + \tau$  then
12:       $b.\text{matched} \leftarrow b.\text{matched} \cup \{i\}$ 
13:       $U \leftarrow U \cup \{i\}$ 
14:      break
15:    end if
16:  end for
17: end for
18:  $R \leftarrow \emptyset$  ▷ Reconstructed regions
19: for all  $b \in B$  do
20:    $I \leftarrow \text{SortAscending}(b.\text{matched})$ 
21:   if  $I = \emptyset$  then
22:      $R \leftarrow R \cup \{\text{Region}(b.\text{label}, b.\text{coords}, \epsilon, \infty)\}$ 
23:     continue
24:   end if
25:    $C \leftarrow \text{SplitIntoContiguousRuns}(I)$ 
26:   for all  $c \in C$  do
27:      $W \leftarrow \{T[k] \mid k \in c\}$ 
28:      $S \leftarrow \text{JoinWordsUsingBreakTypes}(W)$ 
29:      $R \leftarrow R \cup \{\text{Region}(b.\text{label}, b.\text{coords}, S, \min(c))\}$ 
30:   end for
31: end for
32:  $O \leftarrow \{i \mid i \in [1, |T|] \wedge i \notin U\}$ 
33: if  $O \neq \emptyset$  then
34:    $C_o \leftarrow \text{SplitIntoContiguousRuns}(O)$ 
35:   for all  $c \in C_o$  do
36:      $W \leftarrow \{T[k] \mid k \in c\}$ 
37:      $S \leftarrow \text{JoinWordsUsingBreakTypes}(W)$ 
38:      $(\bar{x}, \bar{y}) \leftarrow \text{MeanCentroid}(W)$ 
39:      $R \leftarrow R \cup \{\text{Region}(\text{"Paragraph"}, [\bar{x}, \bar{y}, \bar{x}, \bar{y}], S, \min(c))\}$ 
40:   end for
41: end if
42: Sort  $R$  in ascending order of sort key
43:  $D.\text{annotations} \leftarrow R$ 
44:  $D.\text{wordsFound} \leftarrow |T|$ 
return  $D$ 

```

---