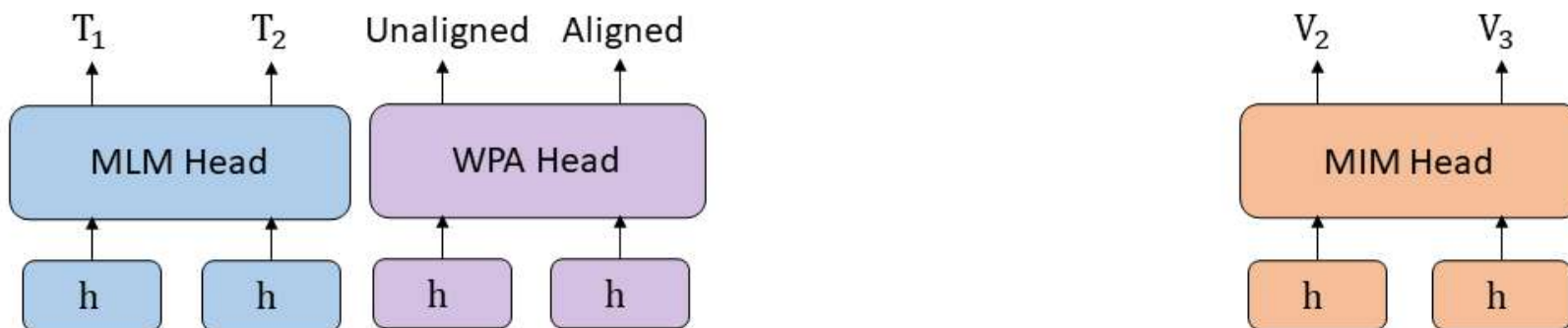


Pre-training Objectives							
-	-	t_3	-	t_5	-	t_7	Masked Visual-Language Modeling
-	r_2	-	r_4	-	-	-	r_8
-	0	-	1	0	-	1	-
							Cross-modal Alignment Identification (0: Mis-aligned, 1: Aligned)

Fine-tuning Tasks								Semantic Entity Recognition (H:Header, Q:Question, A:Answer, O:Other)
<i>O</i>	<i>H</i>	<i>O</i>	<i>Q</i>	<i>A</i>	<i>O</i>	<i>Q</i>	<i>A</i>	Relation Extraction (0=None, 1:Key-Value Pair)
-	-	-	-	-	-	-	-	
-	-	-	-	-	-	-	-	
-	-	-	-	-	-	-	-	
-	0	-	-	-	-	-	-	
-	0	-	1	-	-	-	-	
-	-	-	-	-	-	-	-	
-	0	-	0	0	-	-	-	
-	0	-	0	0	-	1	-	
.....							

Pre-training Objectives

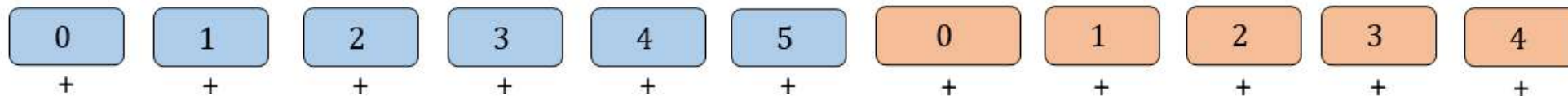


Multimedial Transformer

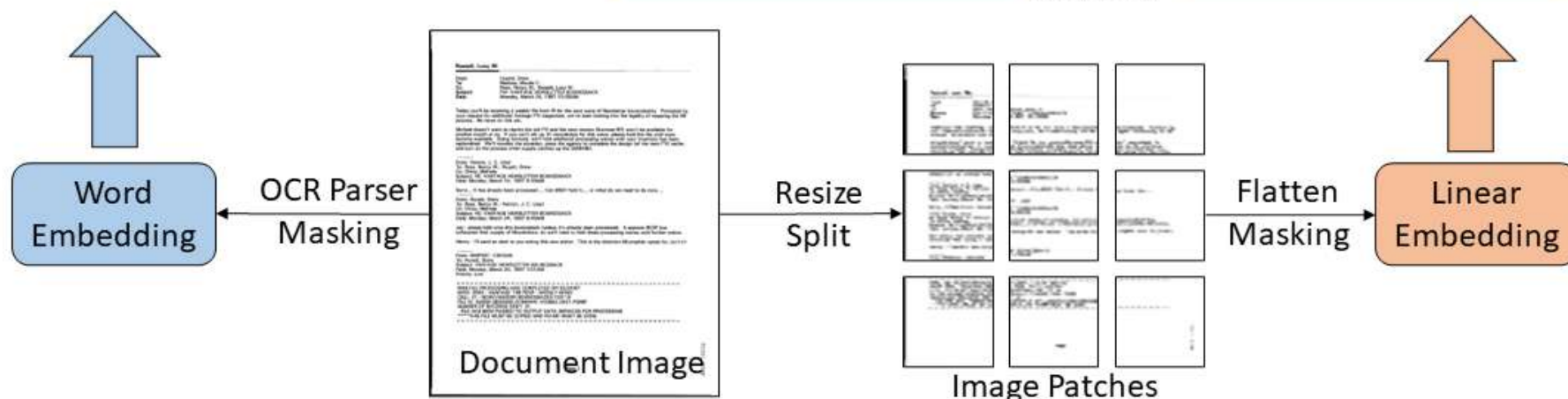
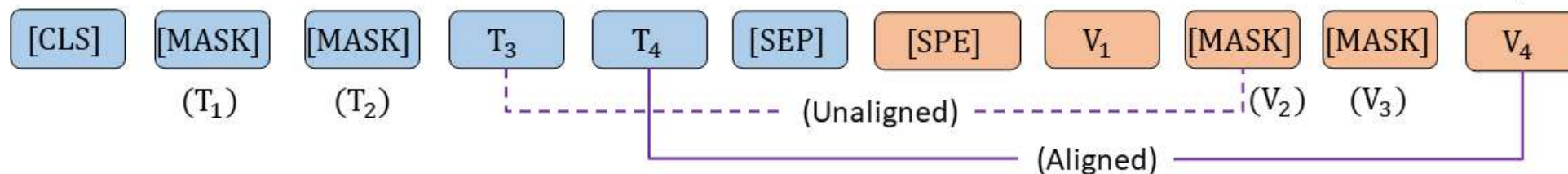
2D Position Embedding



1D Position Embedding



Word/Patch Embedding



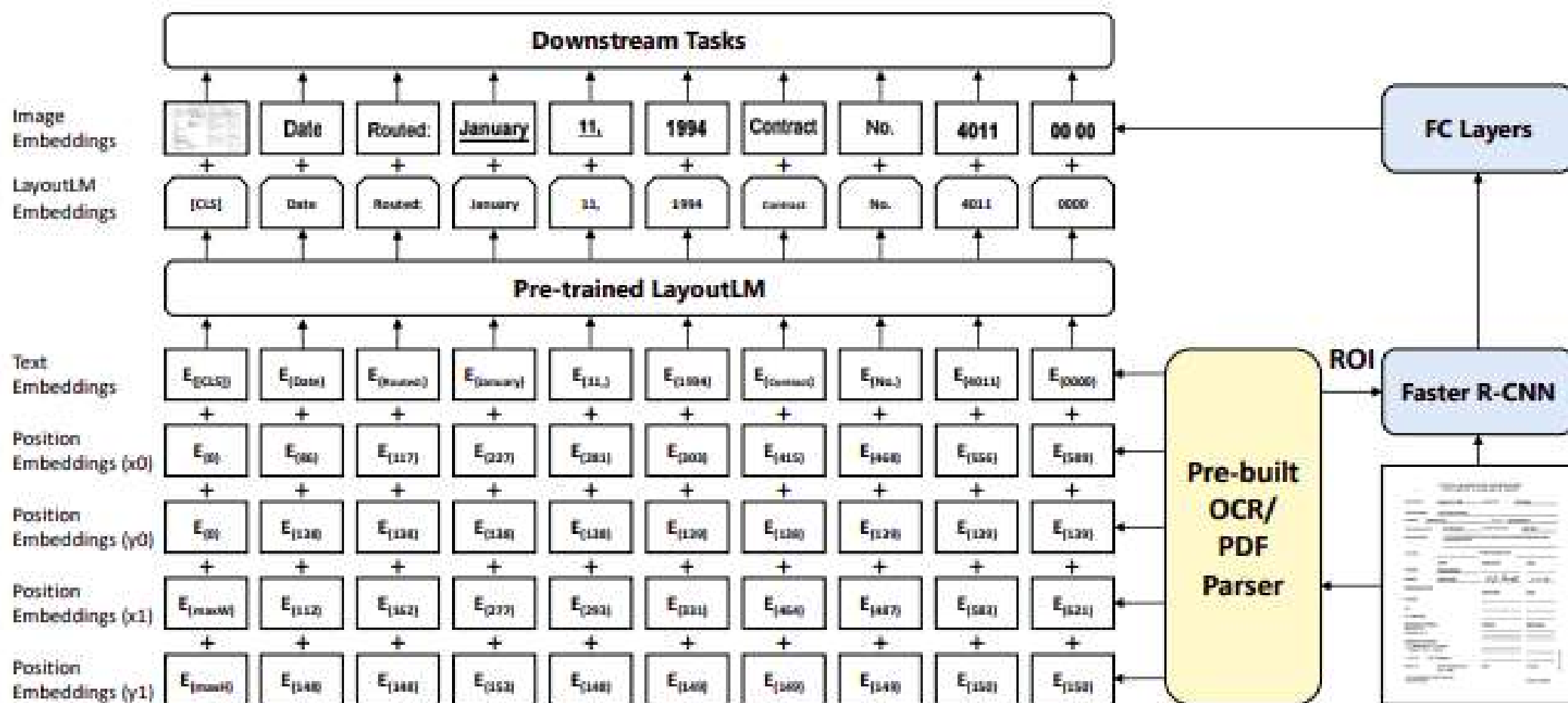


Figure 2: An example of LayoutLM, where 2-D layout and image embeddings are integrated into the original BERT architecture. The LayoutLM embeddings and image embeddings from Faster R-CNN work together for downstream tasks.

BIOES

- BIOES tagging (BIO variants)

- Tagged by one word

- **B** : token that begins a span
 - **I** : tokens inside a span
 - **O** : tokens outside of any span
 - **E** : token that ends a span
 - **S** : token that gets single span

Words	Tag	BIO Tag	BIOES Tag
Jane	PER	B-PER	B-PER
Villanueva	PER	I-PER	E-PER
of		O	O
United	ORG	B-ORG	B-ORG
Airlines	ORG	I-ORG	I-ORG
Holding	ORG	I-ORG	E-ORG
discussed		O	O
the		O	O
Chicago	LOC	B-LOC	S-LOC
route		O	O
.		O	O