



微信扫一扫
关注该公众号



文 | 付奶茶

随着最近几年多模态大火的，越来越多的任务都被推陈出新为多模态版本。譬如，传统对话任务，推出了考虑视觉信息的多模态数据集；事件抽取，也推出视频形式的多模态版本；就连 grammar induction（语法归纳），也有了多模态版的（详见 NAACL'2021 best paper）。

然而，多模态大火是最近的事情，但它并不是近两年才有的什么新技术。如果是想要对这一领域有比较深的研究，甚至想要做出工作、有所创新，那仅仅了解多模态最近两年几个大火的 多模态模型显然是不够的。

事实上，有些任务已经天生就是多模态很多年了。早在多模态成为焦点之前，就已经默默被研究二十来年了。比如，智能文档（Document AI）技术，所谓智能文档技术，也就是自动理解、分析业务文档需求，文档内容可包含文字、图片、视频等多种形式。由于理解多模态形式的多模态形式文的需求其实广泛长期存在，所以智能文档技术很多年来都是几个大厂的 研究重点之一。近年来，深度学习技术的普及也更好地推动了例如文档布局分析、可视化信息提取、文档可视化问答、文档图像分类等智能文档算法的发展。近期，微软亚研院发表了一篇综述，简要回顾了一些有代表性的 DocumentAI 的模型、任务和基准数据集。小编认为这篇综述的总结体系非常扎实，是值得细细阅读的多模态相关综述，故与各位分享。

论文标题：

Document AI: Benchmarks, Models and Applications

论文链接：

<https://arxiv.org/abs/2111.08609>

Document AI 发展历程

作者概述智能文档的发展大致经历了以下三个阶段：

第一阶段：启发式阶段

20 世纪 90 年代初，研究人员主要使用基于规则的启发式 (Heuristic rule-based document layout analysis) 来理解和分析文档，通过手动观察文档的布局信息，从而总结出一些启发式规则。启发式规则方法主要使用固定的布局信息来处理文档，方法较为固定，定制的规则可扩展性较差、通用性较差。

基于启发式规则的文档的布局大致分为三种方式：

(1)自顶向下:文档图像逐步划分到不同的区域,递归执行切割直到该区域被划分为预定义的标准，通常是块或列。例如projection profile,采用X-Y cut算法对文档进行剪切,通常用于文本区域和行距固定的结构化文本,对特定格式的文档进行更快、更有效的分析,但其对边界噪声敏感,对倾斜文本的处理效果不佳。

(2)自底向上:使用像素或组件作为基本单元,将其分组并合并成一个更大的同质区域,自底向上方法虽然需要更多的计算资源,但更通用,可以覆盖更多具有不同布局类型的文档。

(3)混合策略:将自上而下和自下而上相结合，例如Okamoto & Takahashi使用分隔符和空格来切割块，并将内部组件进一步合并到每个块中的文本行中,进而解析文档的布局。

第二阶段:机器学习阶段

直到从 2000 年来 随着机器学习技术的发展，以机器学习模型逐渐成为文档处理的主流方法。研究者设计功能模板以了解不同功能的权重，进而理解和分析文档的内容和布局。

基于机器学习的文档分析过程通常分为两个阶段：

1)对文档图像进行分割，获得多个候选区域；
2)对文档区域进行分类和区分，如文本块和图像。

尽管带注释的数据被用于监督学习，并且以前的方法可以带来一定程度的性能改进，但是由于缺乏定制规则和训练样本数量，通用性仍然不令人满意。此外，不同类型文档的迁移和适应成本相对较高，这使得以前的方法不适合广泛的商业应用。

第三阶段:深度学习阶段

随着深度学习的发展和大量未标注电子文档的积累，可以通过工具 HTML/XML 提取、PDF 解析器、OCR 等提取不同类型的文档中的内容，其文本内容、布局信息和基本图像信息等基本组织良好，然后对大规模深度神经网络进行预训练和微调，以完成各种下游文档 AI 任务，包括文档布局分析、视觉信息提取、文档视觉问答和文档图像分类等。现有的基于深度学习的智能文档模型主要分为两大类：

- 针对特定任务的深度学习模型
- 支持各种下游任务的通用预训练模型

Document AI 的主要任务

Document AI 在我们现实的应用场景主要有以下四类任务：

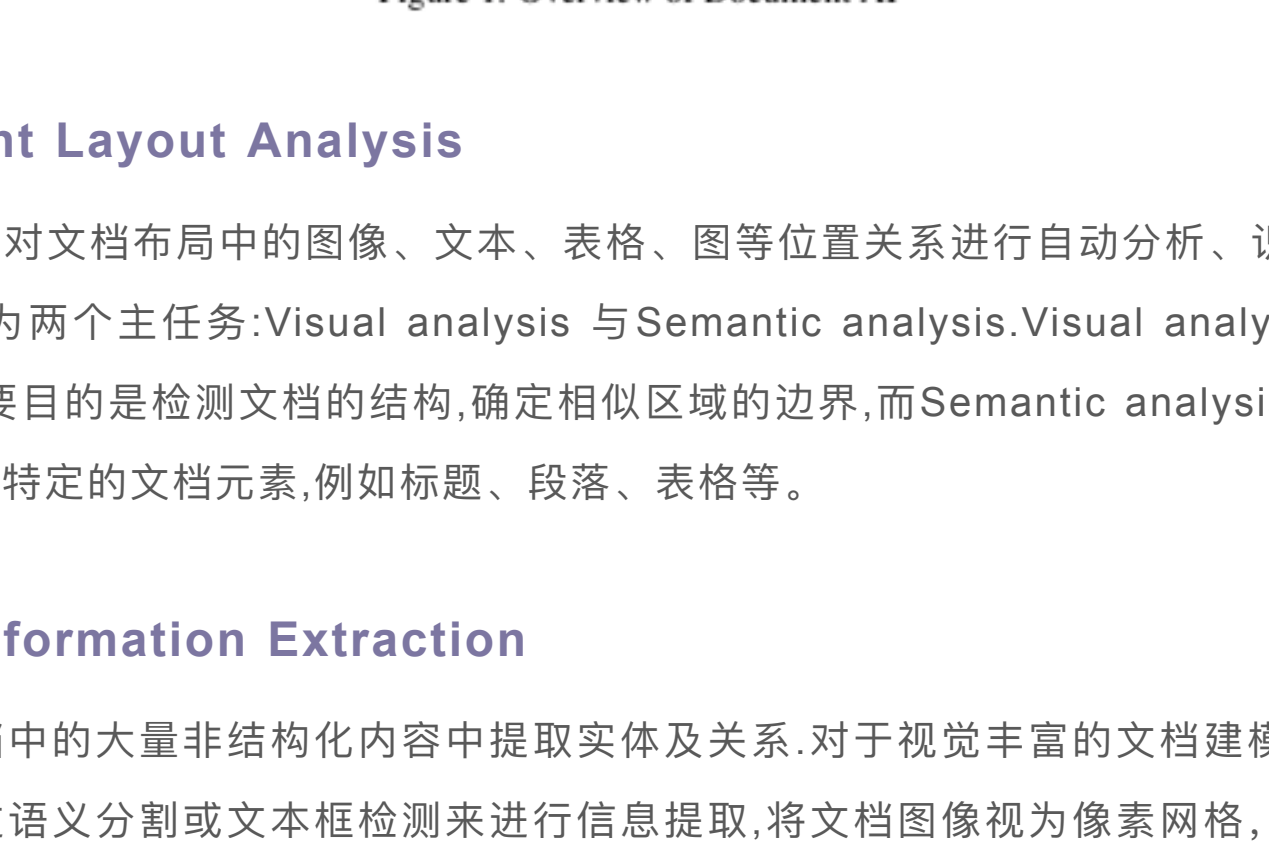


Figure 1: Overview of Document AI

Document Layout Analysis

该任务主要是对文档布局中的图像、文本、表格、图等位置关系进行自动分析、识别、理解的过程。主要分为两个子任务:Visual analysis 与 Semantic analysis. Visual analysis 为视觉元素的分析,主要目的是检测文档的结构,确定相似区域的边界,而 Semantic analysis 为语义分析,检测区域识别特定的文档元素,例如标题、段落、表格等。

Visual Information Extraction

该任务从文档中的大量非结构化内容中提取实体及关系。对于视觉丰富的文档建模为计算机视觉问题,通过语义分割或文本框检测来进行信息提取,将文档图像视为像素网格,将文本特征添加到视觉特征图中。根据文本信息的程度,该任务从字符级发展到单词级,再发展到上下文级。

Document Visual Question Answering

该任务为通过判断识别文本的内部逻辑来回答关于文档的自然语言问题。文档 VQA 中的文本信息在任务中起着至关重要的作用，现有的有代表性的方法都是以文档图像的 OCR 获取的文本作为输入。获得文档文本后，将 VQA 任务建模为不同的问题:主流方法将其建模为机器阅读理解 (MRC) 问题,根据问题从给定文档中提取文本片段作为相应的答案。

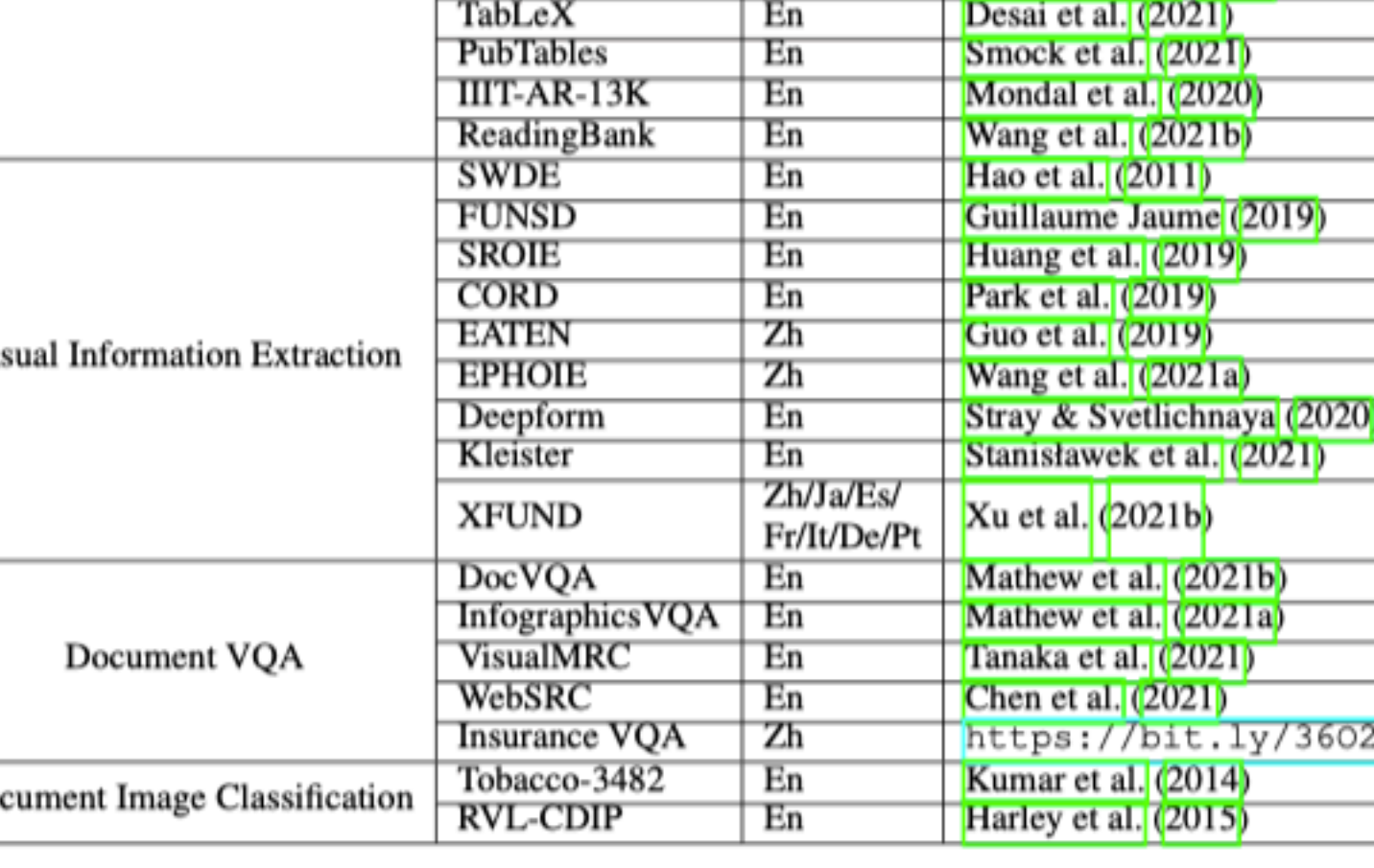


Table 1: Benchmark datasets for document layout analysis, visual information extraction, document visual question answering and document image classification.

Document AI 主流模型

Documents layout analysis with convolutional neural networks

文档布局分析可以看作是对文档图像进行目标检测的任务。将文档中的标题、段落、表格、图表等基本单元是需要检测和识别的对象。Yang 等人将文档布局分析作为像素级的分割任务，利用卷积神经网络进行像素分类，取得了较好的效果。

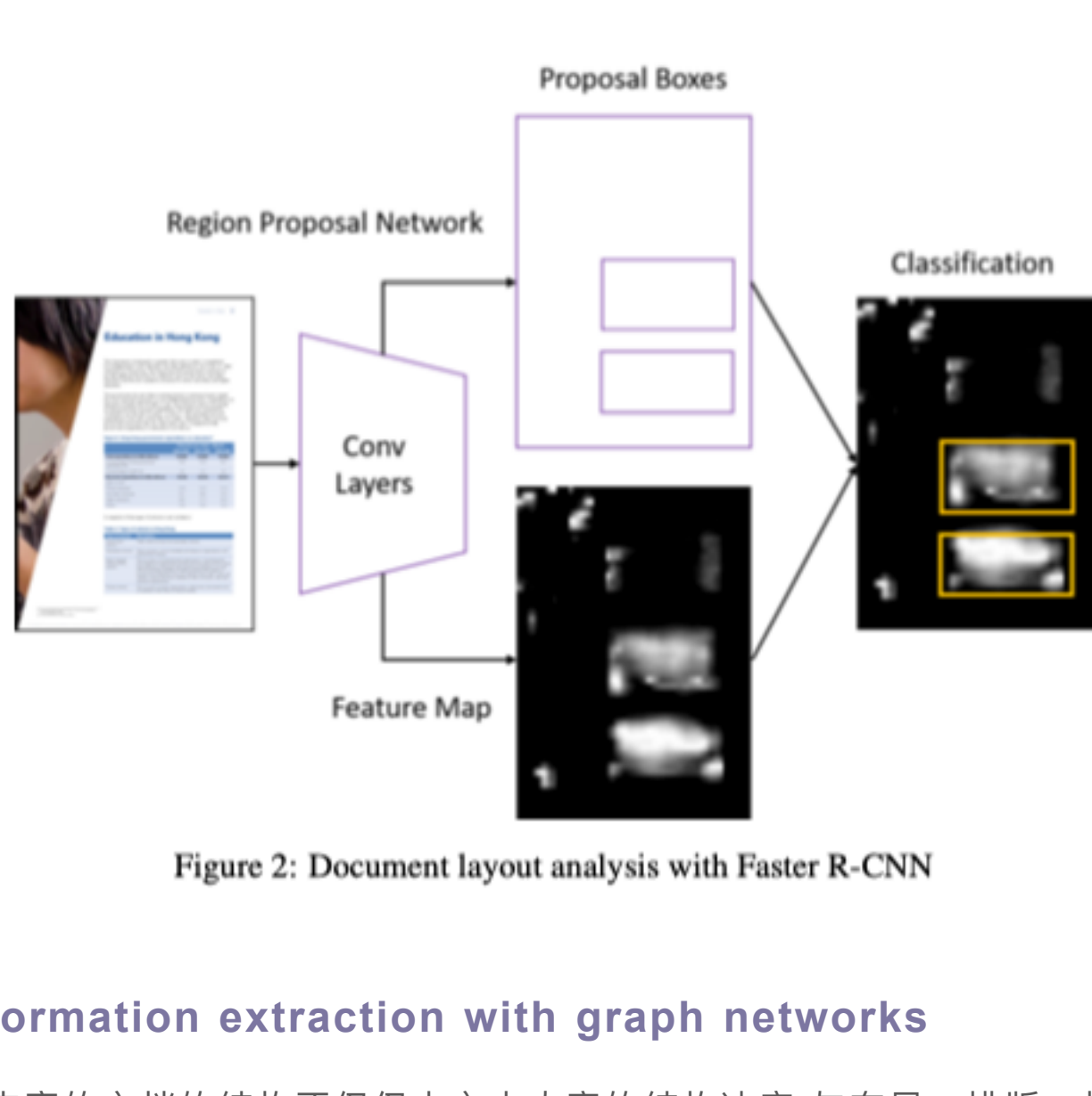


Figure 2: Document layout analysis with Faster R-CNN

Visual information extraction with graph networks

对于视觉信息丰富的文档的结构不仅仅由文本内容的结构决定,与布局、排版、格式、表/图结构等视觉元素同样相关,例如收据、证书、保险文件等.Liu 等人提出的利用图卷积神经网络建模视觉元素丰富的文档,首先通过 OCR 系统获得一组 Text Blocks,每个 Text Block 包含其在图像中与文本内容的坐标信息,将其构成一个完全连通的有向图,即每个 Text Blocks 构成一个节点,通过 Bi-LSTM 获取节点的初始特征,边的初始特征是相邻文本块与当前文本块之间的相对距离以及这两个文本块的长宽比。对“节点-边缘-节点”三元特征集进行卷积,实验表明,视觉信息发挥了主要作用,增加了文本识别相似语义的能力,对视觉信息也起到一定的辅助作用。

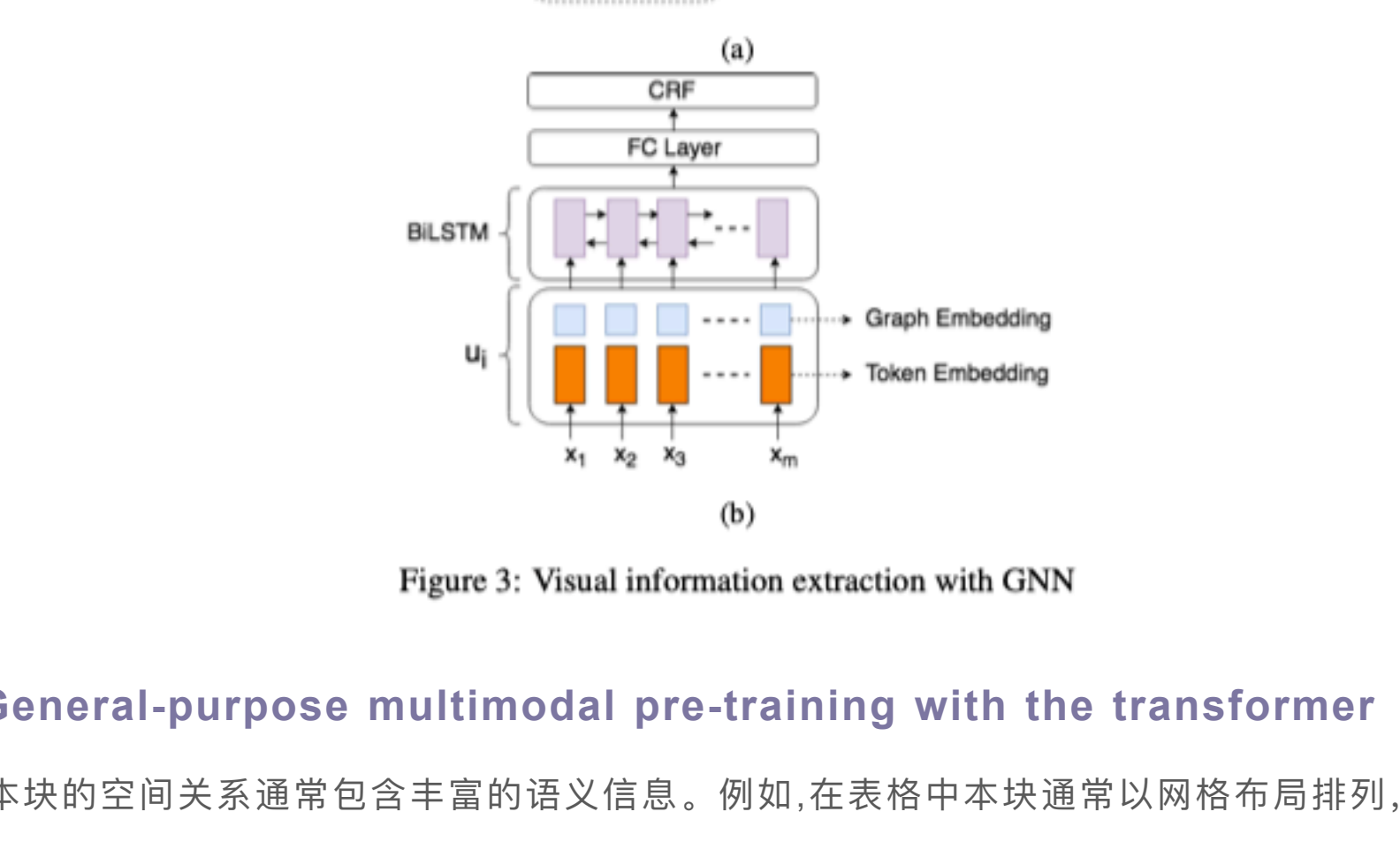


Figure 3: Visual information extraction with GNN

General-purpose multimodal pre-training with the transformer

文本块的空间关系通常包含丰富的语义信息。例如，在表格中本块通常以网格布局排列，标题通常出现在第一列或第一行。不同文档类型之间的布局不变性是通用预训练的一个关键属性。通过预训练与文本自然对齐的位置信息可以为下游任务提供更丰富的语义信息。对于视觉信息丰富的文档，其视觉信息如字体类型、大小、样式等明显的视觉差异，其可以通过视觉编码器提取出来，结合到预训练阶段，从而有效地改善下游任务。为了利用布局和视觉信息，2020 年 Xu 提出通用文档预训练模型 LayoutLM，在已有预训练模型的基础上，增加了 2-D position embedding 和 image embedding。首先根据 OCR 得到的文本边界框得到文本在文档中的坐标。将对应的坐标转换为虚坐标后，模型计算出 x、y、w、h 四个 embedding sublayers 对应的坐标表示，最终的二维位置嵌入是四个子层的 embedding 之和。在 image embedding 中，模型将每个文本对应的边框作为 Faster R-CNN 提取相应的局部特征。特别是，由于 [CLS] 符号用于表示整个文档的语义，因此模型还使用整个文档的 image 作为 image embedding 以保持多模态对齐。Layout 模型在三个下游任务，表理解，票据理解，文档图像分类，都取得了显著的准确率提升。

LayoutLM 的两个自监督预训练任务：Task1:Masked Visual-Language：随机 mask 除了 2D position embedding，以及其他文本的 text embedding，让模型预测 mask 的 token。Task2:Multi-Label Document Classification：给定一组扫描文档的情况下，利用文档标签对训练前的过程进行监督，使模型能够对来自不同领域的知识进行聚类，生成更好的文档级表示。该模型的相关实验表明，利用布局和视觉信息的预训练可以有效地转移到下游任务中。

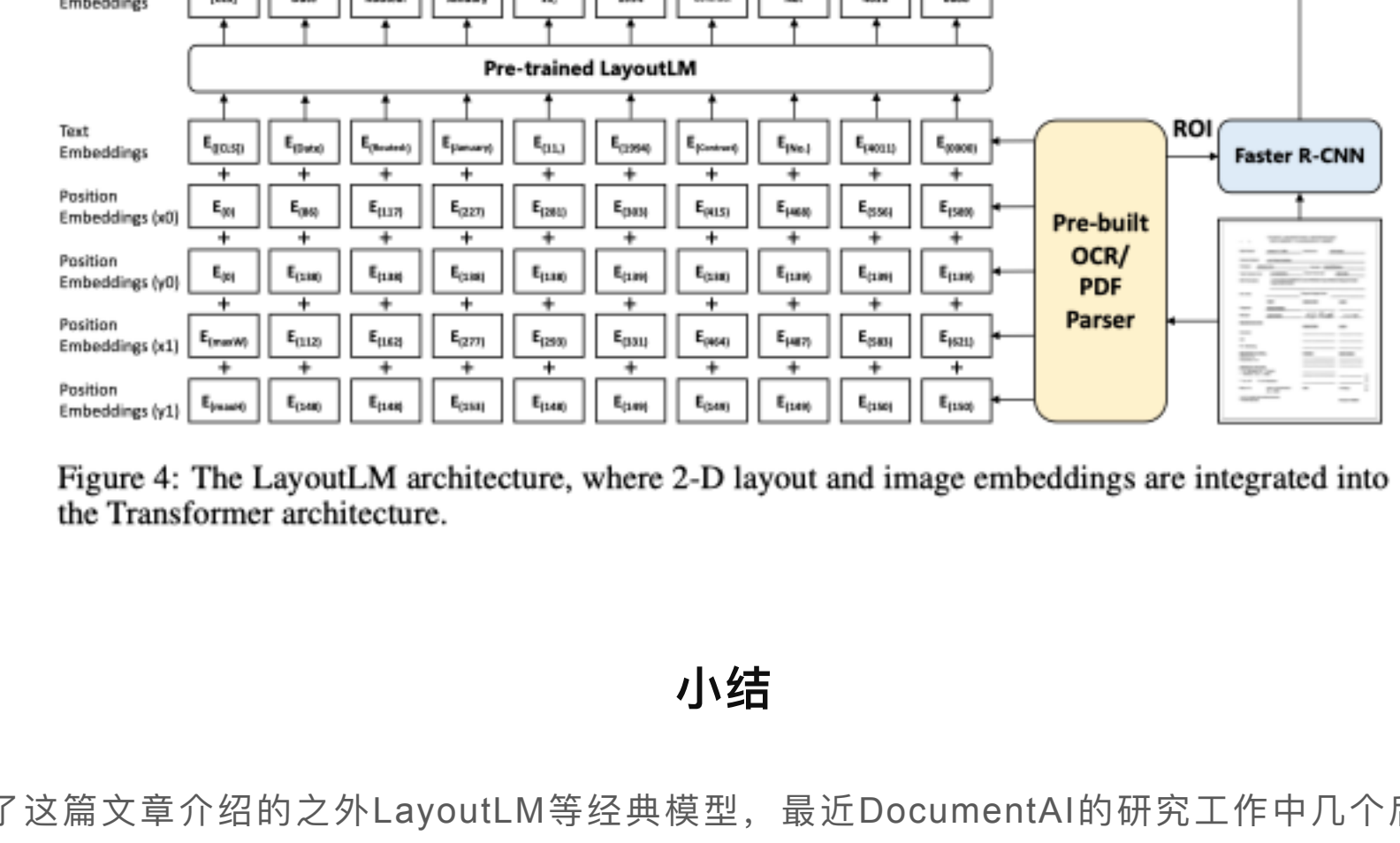


Figure 4: The LayoutLM architecture, where 2-D layout and image embeddings are integrated into the Transformer architecture.

小结

除了这篇文章介绍的之外 LayoutLM 等经典模型，最近 DocumentAI 的研究中几个后起之秀也非常值得关注。例如 LayoutLM 后出现的 LayoutLMv2 以及 LayoutXML，将跨模态对齐的思路贯彻在模型训练的过程中。不仅仅利用文本和布局信息，将图像信息也融合到文档多模态的框架内。除此之外，跨模态文档理解模型 ERINE-Layout，提出阅读顺序预测和细粒度图文匹配两个与训练任务，除了跨模态予以对齐能力外，增加了布局理解能力。我们可以看到，在预训练时代下，DocumentAI 正在逐渐向“多模态文档理解”方向前进，从模态之间的对齐到预测，DocumentAI 将会怎样寻找可以建模的更多元素，挖掘视觉与文本、布局之间的精细关系，变得更加值得期待了。

卖萌屋相关阅读：《别再再搞纯文本了！多模态文档理解要被时代需要！》

萌屋作者：付奶茶

新媒体交叉学科在读 PhD, 卖萌屋十级粉丝修炼上在小编, 目前深耕多模态, 希望可以和大家一起认真科研, 快乐生活!

作品推荐

1. 在斯坦福, 做 Manning 的 phd 要有多难?
2. 史上最大多模态图文数据集发布!

后台回复关键词【入群】

加入卖萌屋 NLP/IR/Rec 与求职讨论群

后台回复关键词【入群】

获取 ACL、CIKM 等各大顶会论文集!

FOLLOW ME

STAR ME

设为星标

推荐给朋友

设置

取消

喜欢此内容的人还喜欢

Nat. Mach. Intell. | MolCLR: 一个用于分子表征学习的自监督框架

DrugAI

《Datawhale 强化学习教程》出版了!

Datawhale