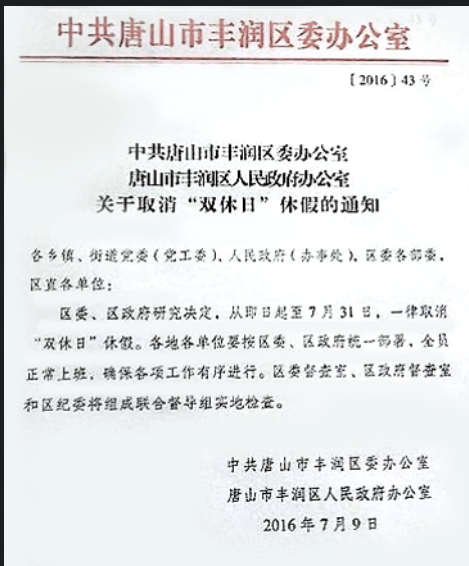


公文要素提取场景

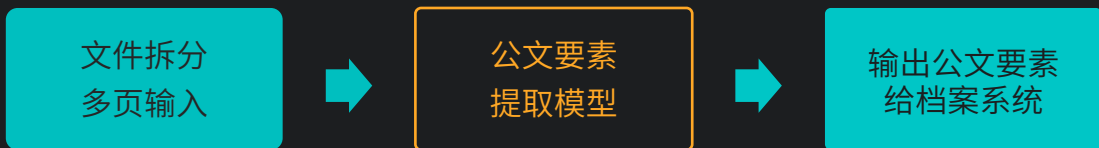
公文是指政府、企业对内、外部的正式发文，有重要信息价值，需要归档并提供全文、关键词检索。



公文要素示例	
抬头	发文单位
文号	发文日期
标题

由于年份和管理机构不同，公文格式也不尽相同，需要模型具备灵活的定位和识别能力

公文要素提取自动化工作流程



公文要素提取建模方案

不限定要素的固定位置，不同版式均可识别：



样本准备

- 1000-1500张图片起，可包含各种版式和无要素页，分布尽可能与实际一致。
- 80%用于训练，20%用于验证
- 准备标注文件（模型学习的正确答案）

训练过程

- 结构化信息提取自动建模工具
- 工具参数配置：Single_value:true; dict_size:300-400; Proba_thresh : 0.35; 后处理步骤output_type:join
- 准备样本后自动训练，耗时约3小时

模型效果

- 准确率约为 **92-97%**