



多模态文字识别课程总结

1. 多模态文字识别概述

- 文字检测 (Text Detection): 找到图像中可能包含文字的区域。
 - 文字识别 (Text Recognition): 将检测到的文字区域转化为可编辑文本。
-

2. 文字检测方法

传统方法

- 基于滑动窗口、边缘检测、连通域分析。
- 缺点：对复杂场景鲁棒性差。

深度学习方法

- 两阶段方法（类似目标检测）：
 - 使用 Faster R-CNN 思路。
 - 单阶段方法：
 - EAST、CTPN、DBNet 等。
 - DBNet 亮点：
 - 基于分割思想，预测每个像素是否为文字。
 - Differentiable Binarization (可微分二值化) 提升检测效果。
-

3. 文字识别方法

- CTC (Connectionist Temporal Classification) :
 - 适合不定长序列建模。
 - Attention 机制：
 - 结合上下文信息识别复杂文字。
 - 常用模型：
 - CRNN、RARE、ASTER。
-

4. 数据与训练

- 常用数据集：
 - ICDAR、SynthText。
 - 训练技巧：
 - 数据增强（旋转、模糊、噪声）。
 - 合成数据生成。
-

5. 多模态融合

- 结合 视觉特征 + 语言模型 (LM):
 - 提升在复杂场景下的识别准确率。
 - 典型做法：
 - 使用预训练语言模型辅助识别结果纠错。
-

6. 实际应用案例

- OCR 系统流程：
 - 图像预处理 → 文本检测 → 文本识别 → 后处理（语言模型纠错）。
- 场景应用：
 - 发票/票据识别、自动驾驶路牌识别、自然场景文字提取。



DBNet 论文总结 (Differentiable Binarization for Scene Text Detection)

1. 背景

- 传统问题：
 - 自然场景文字检测困难：弯曲文字、复杂背景、小字体。
 - 现有方法二值化过程不可导，训练与推理不一致。
- 核心贡献：

- 🌟提出 可微分二值化 (Differentiable Binarization, DB), 端到端优化, 提高检测精度和效率。

2. 🌟 模型结构

1. Backbone

- 轻量级 CNN (ResNet/FPN) 。
- 提取多尺度特征。

2. 预测模块

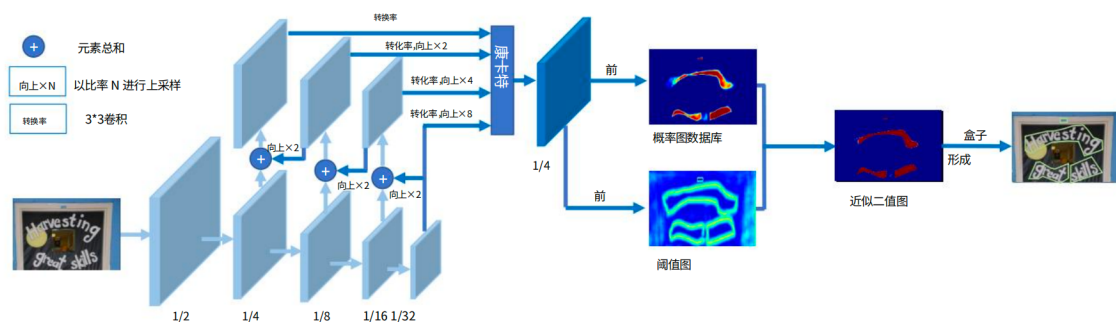
- 预测两个图：
 - 概率图 (P): 像素属于文本的概率。
 - 阈值图 (T): 每个像素的二值化阈值。

3. 可微分二值化函数

- 公式：

$$B(x) = 1 / (1 + \exp(-k * (P(x) - T(x))))$$

- k 控制平滑程度 (越大越接近硬阈值) 。
- 可导 → 允许端到端训练。



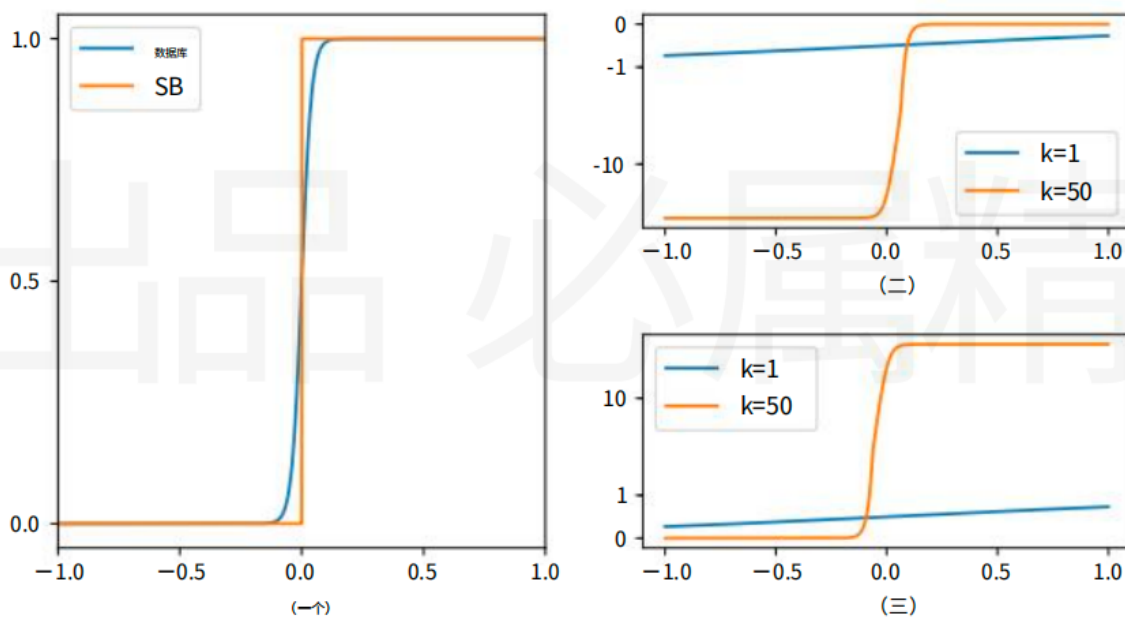


图 4:可微分二值化及其导数。(a)标准二值化 (SB)和可微二值化 (DB)的数值比较。(b) l^+ 的导数。(c) l^- 的导数。

3. 训练与推理

- 损失函数
 - BCE 损失 (概率图)。
 - L1 损失 (阈值图)。
 - IoU 损失 (整体边界质量)。
- 推理
 - 通过 DB 生成二值化图。
 - 使用后处理 (连通域、极小区域过滤) 得到文本区域。

4. 实验与效果

- 数据集: ICDAR 2015、ICDAR 2017、Total-Text、CTW1500。
- 结果:
 - 在弯曲文字、复杂背景下显著优于 EAST、PSENet。
 - 精度和速度兼顾:

- ResNet-18: 55 FPS
- ResNet-50: 26 FPS

5. 技术亮点（记忆口诀）

- 双图预测：P 概率图 + T 阈值图
- 可微分二值化：解决训练/推理不一致
- 端到端优化：检测精度更高
- 轻量高效：实时场景可用

👉 记忆口诀：

“双图预测，端到端；可微二值，检测强；轻量高效，场景广。”

ABINet 论文总结

1. 背景与问题

- 传统 OCR 流程：文字检测 → 文字识别
- 挑战：自然场景文字识别中存在弯曲、模糊、背景复杂、上下文依赖等问题。
- 核心思路：将 视觉信息 (Vision) 和 语言信息 (Language) 有机结合，构建多模态识别模型。

2. 模型结构 (ABINet 框架)

整体由三部分组成：

1. Vision Model (视觉模型)

- 输入图像，提取视觉特征。
- 预测初步的字符序列。

2. Language Model (语言模型)

- 基于 Transformer 结构。
- 通过上下文建模预测字符分布。

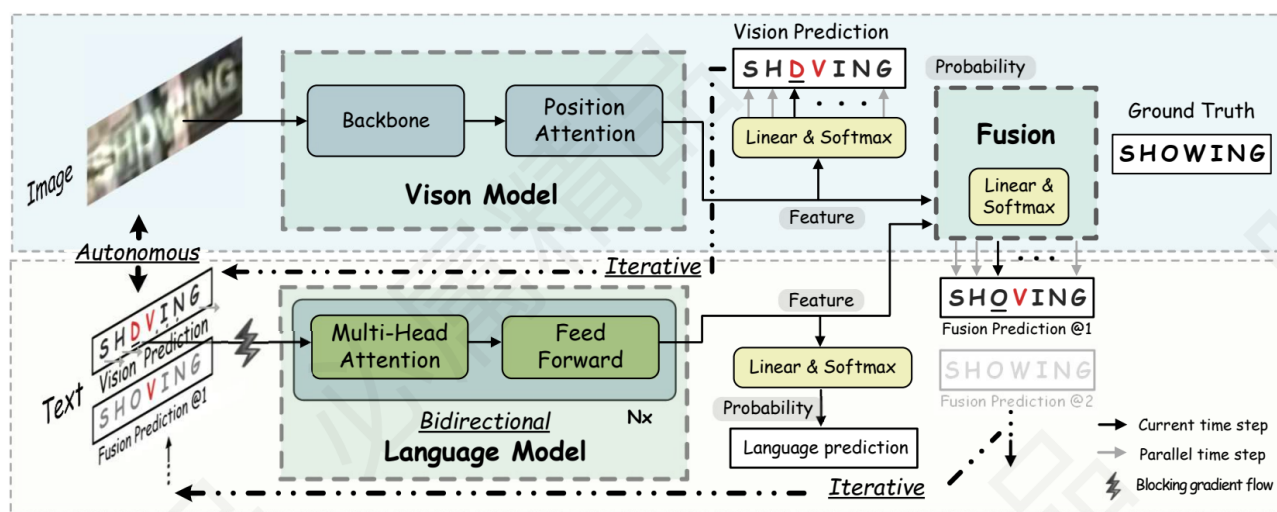
3. Fusion Module (模态融合)

- 将视觉预测与语言预测进行交互。

- 迭代优化，提升识别效果。

关键机制：**Iterative Correction (迭代校正)**

模型多次修正识别结果，使预测更准确。



3. 技术亮点

- 多模态融合：视觉 + 语言的动态交互，而非简单加权。
- 自回归替代：通过并行预测替代传统自回归解码，提高推理速度。
- 训练方式：联合训练，端到端优化。
- 语言先验：内置强大的语言建模能力，能自动纠错。

总结

方面	DBNet (检测)	ABINet (识别)
 任务定位	场景文字 检测 ：找到文字区域	场景文字 识别 ：读出具体的字符
 核心	提出 可微分二值化 (DB) ，解	提出 视觉-语言融合 + 迭代校正 ，提升识别精

贡献	解决训练与推理不一致	度
 模型结构	- Backbone (ResNet/FPN) 提取特征 - 概率图 (P) + 阈值图 (T) - 可微分二值化函数 B(x)	- Vision 模块: 图像特征 → 初步字符序列 - Language 模块: Transformer 上下文建模 - Fusion 模块: 视觉 + 语言交互, 迭代优化
 关键公式/机制	$B(x) = 1 / (1 + \exp(-k * (P(x) - T(x))))$	Iterative Correction (多轮修正预测结果)
 训练方式	- BCE (概率图) - L1 (阈值图) - IoU 损失	- 端到端联合训练 - 并行预测替代自回归
 速度表现	ResNet-18 可达 55 FPS (实时)	非自回归并行预测, 比传统识别更快
 效果表现	弯曲文字/复杂背景检测效果显著提升	多数数据集 SOTA, 复杂场景识别更稳健
 亮点记忆	“双图预测, 端到端, 可微二值, 检测强”	“视觉+语言, 迭代校正, 并行高效, 识别准”



应用场景

文本检测 → OCR 前置模块：
票据、路牌、广告牌

文本识别 → OCR 后置模块：身份证、车牌、
街景文字