# Numerical Formula Recognition from Tables

Qingping Yang[1,2], Yixuan Cao[1,2], Hongwei Li[3], Ping Luos[1,2,4]

[1] Institute of Computing Technology, Chinese Academy of Sciences, [2] University of Chinese Academy of Sciences
[3] Research Department, P.A.I. Ltd., [4] Peng Cheng Laboratory

## Introduction

- Claims over the numerical relationships among some measures are commonly expressed as formulas in tabular forms

- This paper introduces the problem of numerical formula recognition from tables



## Rethinking on Table

- Table is a kind of **language** that adopts a different linguistic paradigm from natural language.

- *Content words* are scattered regularly in table cells, and *visual grammar* express the grammatical relationships among the table cells.



## Challenges

- Recognizing formulas require decoding the visual grammar while simultaneously understanding the textual information.

- Horizontal formulas are common in tables.

- Multiple formulas might appear in the same table cell.

- Formula Complexity



## Methods - TaFor

- Problem Conversion.

  A formula can be defined as:
  $$r = f(e_1, \ldots, e_i, \ldots, e_n)$$
  Converted to a set of triplets as:
  $$\{(r, f^1, e_1), \ldots, (r, f^i, e_i), \ldots, (r, f^i, e_i)\},$$
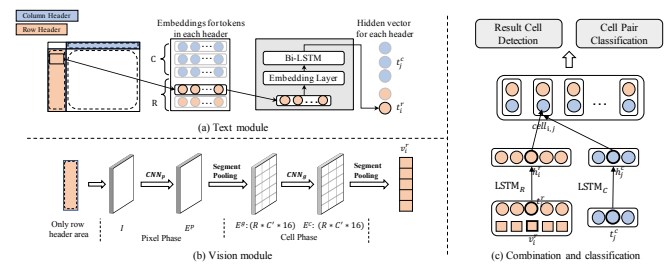  where $r$ is the result cell, $f$ is the formula type, $e$ is the element cell.

- Two Steps.

  1. Result Cell Detection

  2. Cell Pair Classification

  How to encode a table and cell inside it?

- Two-channel Model



(a) Text module

(b) Vision module

(c) Combination and classification

## Experiments

Table 2: Evaluation results.

|  | ± | d | gr | avg | overall |
|---|---|---|---|---|---|
| HHM | 42.57 | 46.29 | 48.78 | 46.37 | 44.08 |
| HSM | 68.00 | 78.97 | 74.45 | 67.12 | 72.05 |
| TaFor | **90.15** | 91.66 | 85.87 | 87.38 | 90.65 |
| HHM + TaFor | 90.02 | **93.58** | **92.19** | **89.18** | **91.31** |

Table 4: Ablation results.

| | Result cell detection | Pair level | Formula level | | | | |
|---|---|---|---|---|---|---|---|
| | | | ± | d | gr | avg | overall |
| TaFor | 96.12 | 95.17 | 90.15 | 91.66 | 85.87 | 87.38 | 90.65 |
| −text | 61.43 | 65.42 | 64.24 | 0 | 0 | 46.40 | 48.78 |
| −vision | 94.42 | 93.93 | 87.86 | 90.89 | 83.69 | 83.59 | 88.77 |

The generalization ability of TaFor