An Integrated Approach of Deep Learning and Symbolic Analysis

for Digital PDF Table Extraction

Mengshi Zhang* University of Texas at Austin Austin, TX, USA mengshi.zhang@utexas.edu Daniel Perelman, Vu Le, Sumit Gulwani Microsoft Redmond, WA, USA {danpere, levu, sumitg}@microsoft.com

*Mengshi Zhang performed this work as part of his internship with the PROSE team at Microsoft. Now, he is a research scientist at Facebook.



Evaluation

Note we do **not** use "intersection-over-union" in order to avoid giving partial credit to tables missing important information, but instead a table is considered correct only on an exact match of all text in the expected table.

Algorithm	Precision	Recall	F ₁
Symbolic	0.315	0.418	0.359
DeepDeSRT (state-of-the-art)	0.178	0.120	0.144
ntegrated (symbolic+our DL)	0.459	0.390	0.422