

The Research Group Mathematics & Data Science

has the honor to invite you to the public defence of the PhD thesis of

Willy Carlos Tchuitcheu

to obtain the degree of Doctor of Sciences

Title of the PhD thesis:

Representation Learning for Table Understanding in Intelligent Document Processing

Supervisor:

Prof. dr. Ann Dooms (VUB) Co-supervisor:

Prof. dr. Tan Lu (VUB)

The defence will take place on

Friday, December 5, 2025 at 3 p.m.

VUB Etterbeek campus, Pleinlaan 2, Elsene, in Learning Theatre (0.04), Learning & Innovation Center (LIC)

The defence can be followed through a live stream: send an email to willy.carlos.tchuitcheu@vub.be to get the link.

Members of the jury

Prof. dr. Jan De Beule (VUB, chair)

Prof. dr. Ann Nowé (VUB)

Prof. dr. Paul Van Eecke (VUB)

Prof. dr. Lynn Houthuys (VUB)

Prof. dr. Apostolos Antonacopoulos (University of Salford, UK)

Dr. Jeroen Jordens (Flanders Make)

Curriculum vitae

Willy Carlos Tchuitcheu obtained his Master's degree in Mathematical Sciences from the African Institute for Mathematical Sciences (AIMS) Rwanda in 2019. Between 2015 and 2020, he gained three years of experience working as a Research Engineer at Camertronix, a research center in Cameroon. In 2021, he began his PhD in the Department of Mathematics and Data Science at VUB. His doctoral work has resulted in three firstauthor articles in international journals (one already published), one patent application, and a Best Poster Award at the 2021 Flanders AI Research (FAIR) Day. In addition, Willy has authored two further journal articles unrelated to this thesis, including one as first author.

Abstract of the PhD research

The volume of business and scientific documents is increasing rapidly, leading to a growing demand for Intelligent Document Processing (IDP) methods. These documents contain a mix of data types, including text, tables, and figures, each with its own structure and semantics. Tables, in particular, differ from text: they are two-dimensional, their meaning does not follow grammatical rules, and they are often permutation-invariant, since swapping rows or columns can preserve their interpretation. A large share of critical information in documents appears in tables (e.g., spreadsheets, reports, invoices, datasheets, and research papers), which humans naturally understand by linking each cell to its corresponding header. We coined the term Table Understanding (TU) principle for this human table-reading ability.

The success of LLMs in understanding text largely relies on rich vector representations called word embeddings. However, meaningful table embeddings are underdeveloped. In practice, LLMs linearize tables into long strings of text, neglecting their structure. This approach loses layout and header information, causing performance to plateau, especially since tables require properties such as permutation invariance, which contrast with the order-sensitive nature of the text that LLMs are trained on.

To address this, we introduce a structure-aware representation learning method for tables. We first model each table as a Heterogeneous Graph Table (HGT) that preserves its two-dimensional layout, cell types, and header relationships. Building on the HGT, we learn context-aware cell embeddings that explicitly bind cells to their headers, aligning with the TU principle. Additionally, we transfer linguistic knowledge from LLMs into these embeddings without flattening the table, thereby combining surface text with structure. Finally, we explore permutation invariance systematically, shuffling rows and columns and assessing the stability of the table embeddings.

Our method demonstrates competitive performance across two main categories of TU tasks: Column Type Annotation (CTA), which falls under the semantic category, and Table Question Answering (TQA), which falls under the reasoning category. Additionally, our approach reduces complexity and enhances robustness against table permutations. These findings indicate that using structure-aware table embeddings grounded in the TU principle can provide a practical solution for reliable, low-latency information extraction in the broader field of IDP.