# NAMED ENTITY RECOGNITION AND RELATION EXTRACTION WITH GRAPH NEURAL NETWORKS IN SEMI STRUCTURED DOCUMENTS

Manuel Carbonell, Pau Riba, Mauricio Villegas,
Alicia Fornés and Josep Lladós

Computer Vision Center Barcelona
omni:us

# Information extraction from semi-structured documents



GOAL: extract information from a document in a structured manner taking layout and semantics in account

# Problem formulation

Information extraction can be reformulated as:

- **Word grouping:** aggregate words into entities
- **Entity labeling:** classify entities into categories (e.g. questions, answers and headers)
- **Entity linking:** find relationships between entities (possibly hierarchical)
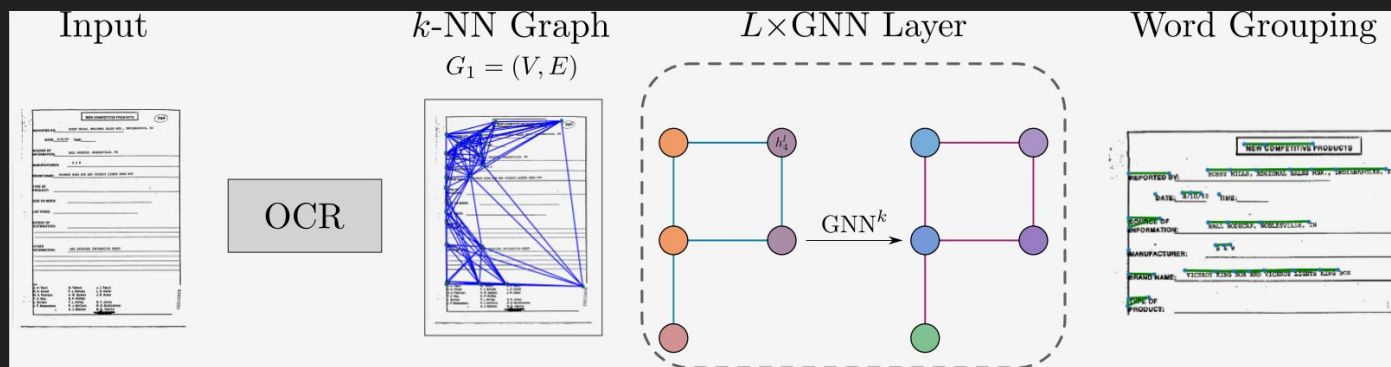
# Examples



In funsd entities are groups of words that denote keys and values to be linked.

In IEHHR entities are groups of words referring to a particular person in a marriage record. Links denote direct relationships between these persons e.g. wife - husband

# Methodology: Word graph

Node features := text box [ *x, y, w, h, word_embed* ]
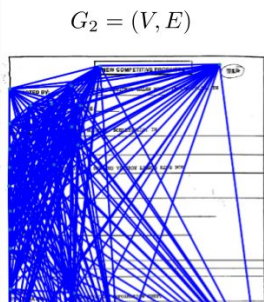


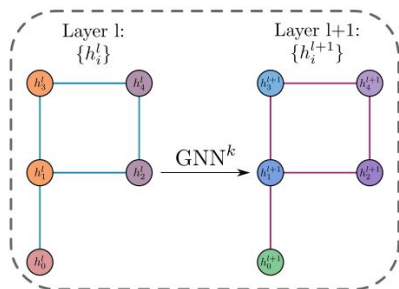GNN is trained for edge classification, to form word groups as connected components

# Methodology: Entity graph

Node features := text box [ *x, y, w, h, entity_embed* ]



GNN is trained for classification of nodes (entity labeling) and edges (entity linking)

# Results

**TABLE I**
RESULTS FOR THE THREE DOCUMENT UNDERSTANDING TASKS ON FUNSD AND IEHHR DATASETS.

| | Word Grouping (ARI) | Entity Labeling (F1) | Entity Linking (F1) | External data | # Params |
|---|---|---|---|---|---|
| **FUNSD** [21] | | | | | |
| [21] | 0.41 | 0.57 | 0.04 | ✓ | 340M |
| [17] | - | 0.79[2] | - | ✓ | 160M |
| **Ours** | 0.65 | 0.64 | 0.39 | - | 201M |
| **IEHHR** [22] | | | | | |
| **Ours** | 0.65 | 0.53 | 0.67 | - | 201M |

# Results

# Conclusions & future work

- GNN node and edge classification provides a promising method for entity recognition and relation extraction in semi structured documents
- The proposed method has been designed for administrative documents but it can also be applied in other domains such as historical manuscripts
- We believe that the obtained results have room for improvement and are limited due to the reduced size of the open available datasets for this type of task.
- Further research is required on a more larger openly available dataset for relation extraction and entity recognition in documents where semantic and spatial information plays a relevant role.