

Graph neural networks for information extraction

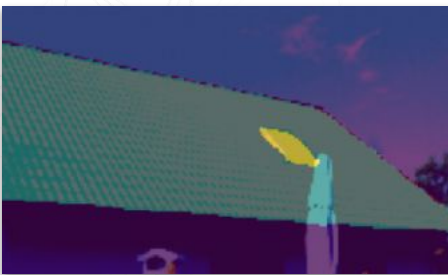
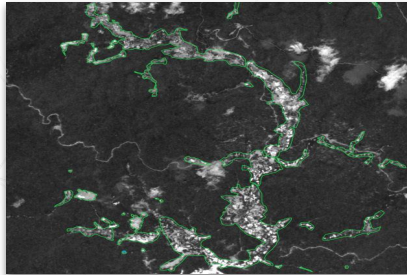
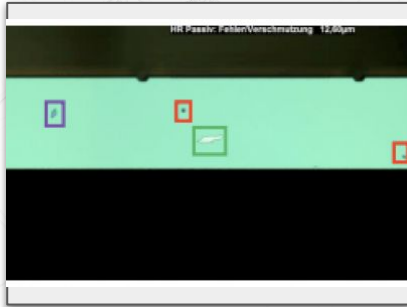
EuroPython 2021

Augusto Stoffel

dida Machine Learning



Deep learning landscape



Computer Vision

[Example Project](#)



Paragraf 9 [7846; 10313] (Schoenheitsreparaturklausel)

§ 9 Instandhaltung der Mieträume und Schönheitsreparaturen

5 ÜBERTRAGUNGSKLAUSEL

6 UMFANGSKLAUSEL

Die Durchführung der Schönheitsreparaturen obliegt dem Mieter. Diese umfassen insbesondere das Tapezieren, Anstreichen der Wände und Decken, das Pflegen der Fußböden, das Streichen der Innentüren, der Fenster und Außentüren...

7 FRISTENKLAUSEL

Bei normaler Benutzung sind die Schönheitsreparaturen, ab Vertragsbeginn gerechnet, in Küche, Bad und WC alle 3 Jahre, für alle übrigen Räume alle 5 Jahre, auszuführen.

8 ENDRENOVIERUNGSKLAUSEL

[...] Demzufolge ist die Mietsache bei Beendigung des Mietverhältnisses unabhängig von der Mietdauer und unabhängig davon, wann zuletzt die vertragsgemäßen Schönheitsreparaturen stattgefunden haben, mit fachmännisch frisch weiß gestrichenen Decken und Wänden sowie im übrigen schadensfrei und gereinigt zurückzugeben.

§ 18 Rückgabe der Mieträume

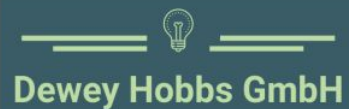
9 FRISTENKLAUSEL

Bei Beendigung des Mietverhältnisses bzw. bei vorherigem Auszug hat der Mieter die Mieträume geräumt, in sauberem Zustand und mit allen - auch den von ihm selbst beschafften - Schlüsseln zurückzugeben. Die Mieträume sind in mangelfreiem Zustand,

Natural Language Processing (NLP)

[Example Project](#)

Sample GNN use case



Dewey Hobbs – Hobbstraße 21 – 12345 Hobbshausen

ABC Chemicals AG
Kochstraße 2
5678 Salzstadt

Dewey Hobbs GmbH
Deweystraße 21
12345 Hobbshausen

Tel.: 0211 12345 67
E-Mail: accounting@hobbs.de
Internet: www.deweyhobbs.de

Rechnung

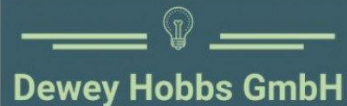
Rechnung Nr. 2020-06-1001 Bestellung-Nr.: KA12345 Kunden-Nr.: 1001

Bitte bei Zahlungen und Schriftverkehr angeben!

Datum: 07.06.2020

Pos	Leistung	MwSt.	Einzelpreis	Anzahl	Gesamtpreis
1	Röntgenblitzröhre ZW 367 H 01G 12 120kV /28mm / 20 cm	19 %	11,00 EUR	2	22,00 EUR
2	50kg Seesand reinst	19 %	22,00 EUR	1	22,00 EUR
3	Salzsäure 1 mol / 10l , 20l Standardlösung 7647-01-0	19 %	33,00 EUR	2	66,00 EUR

Sample GNN use case



Dewey Hobbs GmbH

Dewey Hobbs – Hobbssstraße 21 – 12345 Hobbshausen

ABC Chemicals AG
Kochstraße 2
5678 Salzstadt

Dewey Hobbs GmbH

Deweystraße 21
12345 Hobbshausen

Tel.: 0211 12345 67
E-Mail: accounting@hobbs.de
Internet: www.deweyhobbs.de

Rechnung

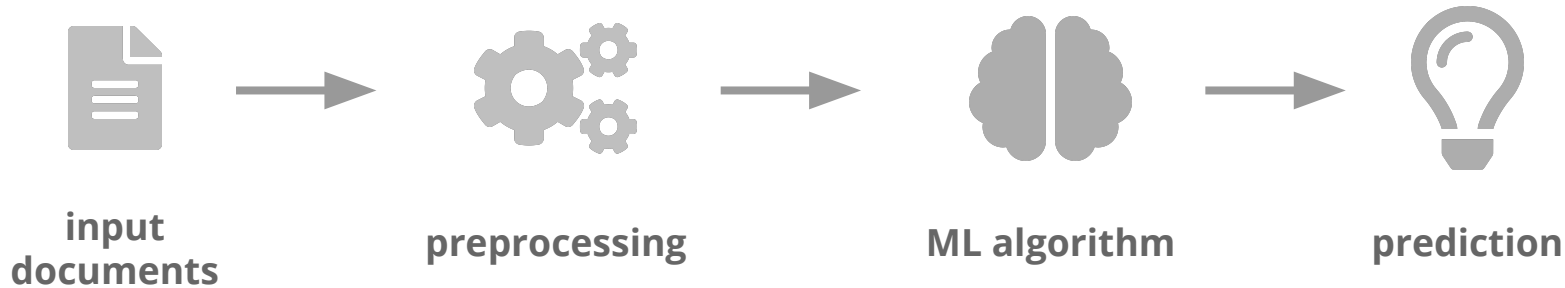
Rechnung Nr. 2020-06-1001 Bestellung-Nr. KA12345 Kunden-Nr.: 1001

Bitte bei Zahlungen und Schriftverkehr angeben!

Datum: 07.06.2020

Pos	Leistung	MwSt.	Einzelpreis	Anzahl	Gesamtpreis
1	Röntgenblitzrohre ZAV 367 H 01G 12 120kV/28mm/20 cm	19 %	11,00 EUR	2	22,00 EUR
2	50kg Seesand reines	19 %	22,00 EUR	1	22,00 EUR
3	Salzsaure 1 mol/l 10l 20 Standardlösung 7647-01-3	19 %	33,00 EUR	2	66,00 EUR

A typical information extraction pipeline



Where do graphs enter in this picture?

- Encode information in a graph
- Use graph algorithms to process the information

A typical information extraction pipeline



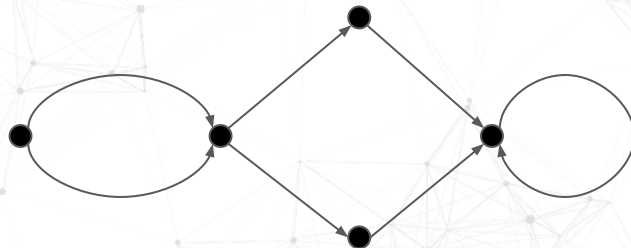
Where do graphs enter in this picture?

- Encode information in a graph
- Use graph algorithms to process the information

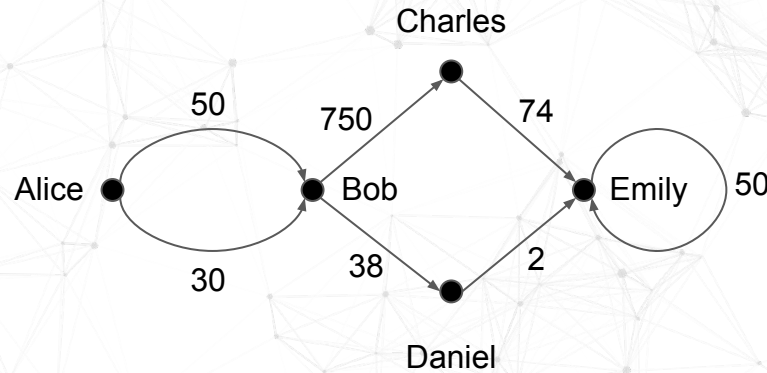
Agenda:

- Graphs and GNNs
- Comparison with convolutional networks
- Use case

What is a graph?



What is a graph?

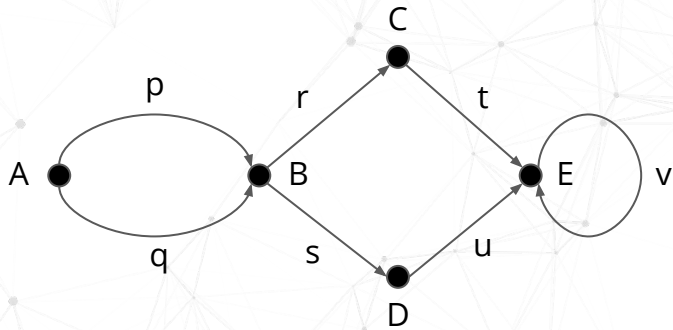


Variations: simple graph, multigraph, pseudograph; directed or not

What is a graph?

Definition: A directed pseudograph, or simply graph, is

- a set \mathcal{V} of vertices
- a set \mathcal{E} of edges
- functions source, target: $\mathcal{E} \rightarrow \mathcal{V}$



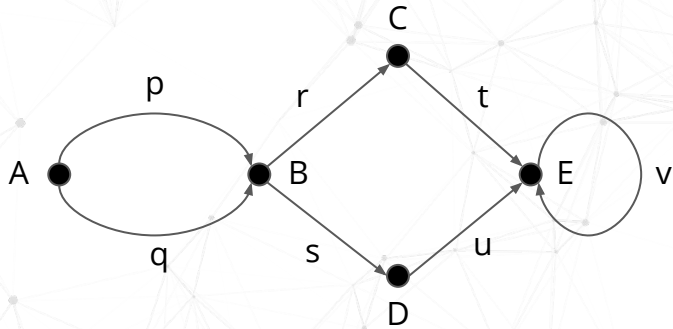
What is a graph?

Definition: A directed pseudograph, or simply graph, is

- a set \mathcal{V} of vertices
- a set \mathcal{E} of edges
- functions source, target: $\mathcal{E} \rightarrow \mathcal{V}$

Example:

- $\mathcal{V} = \{A, B, C, \dots, E\}$
- $\mathcal{E} = \{p, q, r, \dots, v\}$
- source(p) = A
- target(p) = B
- source(v) = E
- target(v) = E



What is a graph?

Definition: A directed pseudograph, or simply graph, is

- a set \mathcal{V} of vertices
- a set \mathcal{E} of edges
- functions source, target: $\mathcal{E} \rightarrow \mathcal{V}$

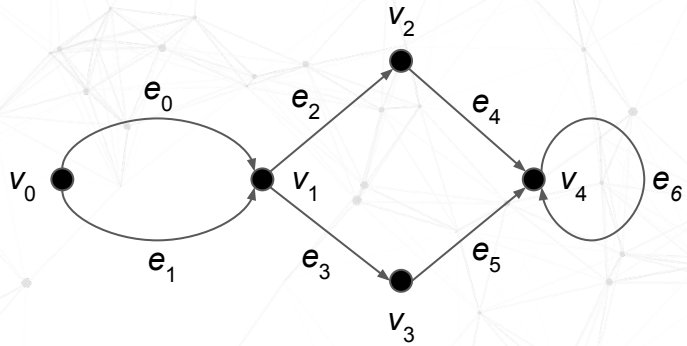
Example:

- $\mathcal{V} = \{A, B, C, \dots, E\}$
- $\mathcal{E} = \{p, q, r, \dots, v\}$
- $\text{source}(p) = A$
- $\text{target}(p) = B$
- $\text{source}(v) = E$
- $\text{target}(v) = E$

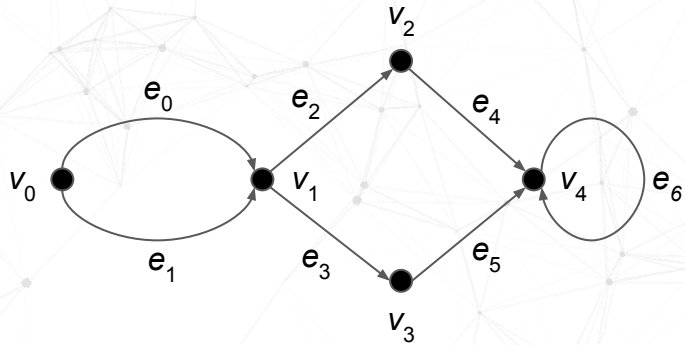
Note: Often there are labels associated to vertices and edges.

- $L_0: \mathcal{V} \rightarrow \{\text{vertex labels}\}$
- $L_1: \mathcal{E} \rightarrow \{\text{edge labels}\}$

Representing graphs in Python

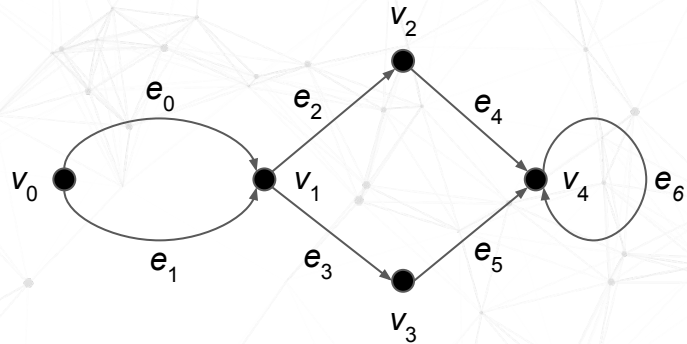


Representing graphs in Python



```
adjacencies = [  
    (0, 1), # e_0  
    (0, 1), # e_1  
    (1, 2), # e_2  
    ...,  
    (4, 4), # e_6  
]
```

Representing graphs in Python



A graph with nodes labeled by NodeT and edges labeled by EdgeT can be modeled as follows:

```
class Graph(Generic[NodeT, EdgeT]):
```

```
    nodes: List[NodeT]
```

```
    edges: List[EdgeT]
```

```
    adjacencies: List[Tuple[int, int]]
```

```
    def source(self, k: int) -> int:
        """Source of the kth edge"""
        return self.adjacencies[k][0]
```

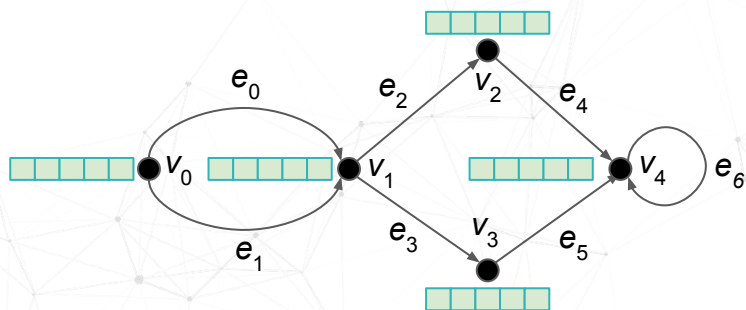
```
    def target(self, k: int) -> int:
        """Target of the kth edge"""
        return self.adjacencies[k][1]
```

```
adjacencies = [
    (0, 1), # e_0
    (0, 1), # e_1
    (1, 2), # e_2
    ...,
    (4, 4), # e_6
]
```

Graph neural networks

Let's consider the *node classification* problem:

- Input: a graph G of type `Graph[Tensor, EnumT]`
- Output: a labeling of its vertices, $L: \mathcal{V} \rightarrow \{0, 1\}$

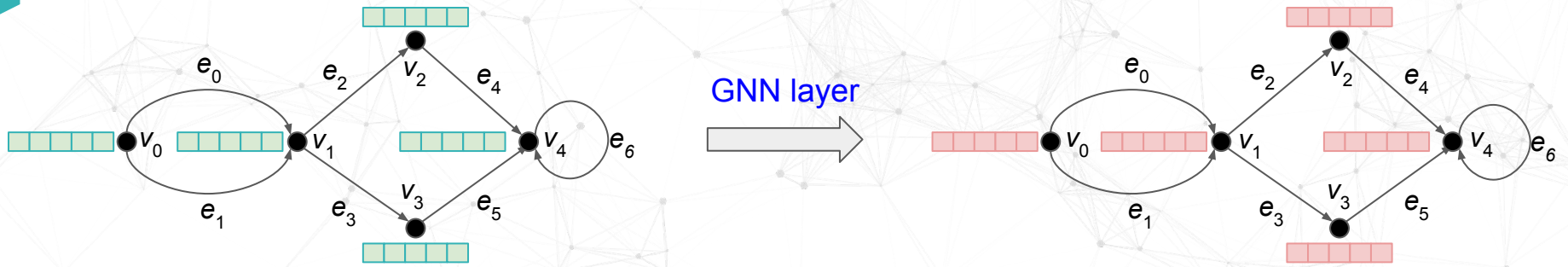


GNN



```
labels = [  
    0, # v_0  
    1, # v_1  
    0, # v_2  
    0, # v_3  
    1  # v_4  
]
```

Graph neural networks



For instance: $y_v = \sigma \left(\sum_{e: w \rightarrow v} \frac{1}{\text{deg}_t^+(v)} W_{t(e)} x_w \right)$ where

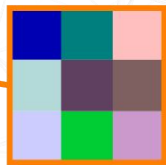
- y_v is the output feature vector of node v
- x_w is the input feature vector of node w
- $t(e)$ is the type of the edge $e: v \rightarrow w$
- W_t is a learned weight matrix corresponding to edge type t
- $\text{deg}_t^+(v)$ is the number of incoming edges of type t for node v
- σ is an activation function

Comparison with convolutional networks

Graph (convolutional) networks generalize CNNs from computer vision



Original image



A 3×3 patch

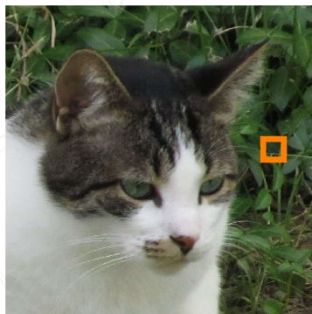
$$y_{i,j} = \sigma \left(\sum_{\substack{k=-1,0,1 \\ l=-1,0,1}} W_{k,l} x_{i+k,j+l} \right)$$

where

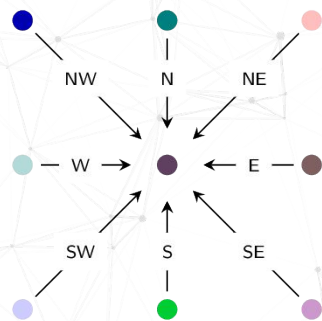
- $y_{i,j}$ is the output feature vector of pixel (i, j)
- $x_{i,j}$ is the input feature vector of pixel (i, j)
- $W_{k,l}$ are learned weight matrices
- σ is an activation function

Comparison with convolutional networks

Graph (convolutional) networks generalize CNNs from computer vision



Original image



A graph representing the 3×3 patch

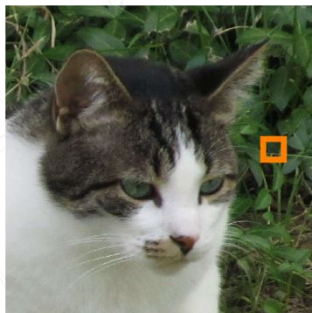
$$y_{i,j} = \sigma \left(\sum_{\substack{k=-1,0,1 \\ l=-1,0,1}} W_{k,l} x_{i+k,j+l} \right)$$

where

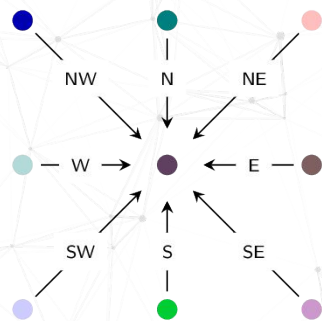
- $y_{i,j}$ is the output feature vector of pixel (i, j)
- $x_{i,j}$ is the input feature vector of pixel (i, j)
- $W_{k,l}$ are learned weight matrices
- σ is an activation function

Comparison with convolutional networks

Graph (convolutional) networks generalize CNNs from computer vision



Original image



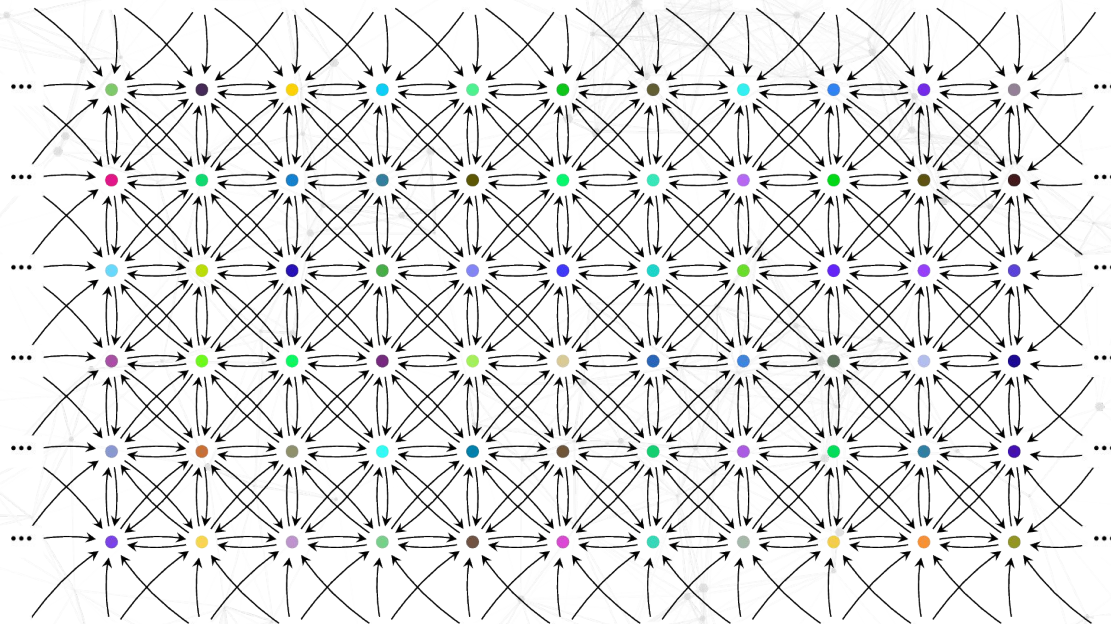
A graph representing the 3 × 3 patch

$$y_p = \sigma \left(W_O x_p + \sum_{q \rightarrow p} W_t x_q \right)$$

where

- x_p is the input feature vector of pixel p
- y_p is the output feature vector of pixel p
- $t \in \{N, NW, W, \dots\}$ is a cardinal direction
- W_t are learned weight matrices
- σ is an activation function

Comparison with convolutional networks

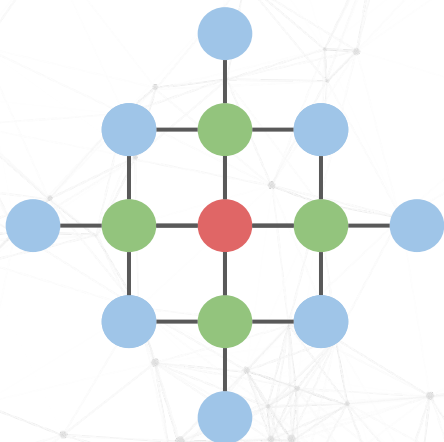


A graph representing the entire image

Use case: information extraction from tables

Idea: Model the document structure as a graph

- each word is a node in the graph
- neighbouring words are connected by edges



Dewey Hobbs GmbH

Dewey Hobbs – Hobbsstraße 21 – 12345 Hobbshausen
ABC Chemicals AG
Kochstraße 2
5678 Salzstadt

Dewey Hobbs GmbH
Deweystraße 21
12345 Hobbshausen
Tel.: 0211 12345 67
E-Mail: accounting@hobbs.de
Internet: www.deweyhobbs.de

Rechnung

Rechnung-Nr.: 2020-06-1001 Bestellung-Nr.: KA12345 Kunden-Nr.: 1001
Bitte bei Zahlungen und Schriftverkehr angeben! Datum: 07.06.2020

Pos	Leistung	MwSt.	Einzelpreis	Anzahl	Gesamtpreis
1	Röntgenblitzröhre ZW 367 H 01G 12 (20kV / 28mm / 20 cm)	19 %	11001,80 €	2	22,00 EUR
2	50kg Seesand reines	19 %	22001,80 €	1	22,00 EUR
3	Salzsäure 1 mol/l 10l 20l Standardlösung 7647-01-0	19 %	88001,80 €	1	66,00 EUR

Use case: information extraction from tables

Nodes:

- Each word is a node in the graph

ACME US&A
999 Supreme Industrial Road , Anderson , South
Carolina .
12345 .UNITED STATES
Tel: +1-123-234-3456 Fed ID 11-123456
Fax: Tax Exempt 123345-6789
VAT Reg. No.:

Seller
We-Supply
NC Suppliers Corporation
P.O. Box 123456
Atlanta
GA 12345-98765
UNITED STATES
Phone: 111-222-3334
Fax:

Ship Via
Road

Terms Of Delivery
Free on board

Order

Our Order Date	Purchase Order
1/2/20	A12345
Revision	
1	

Supplier No
WESUPPLY

Delivery Address
ACME US&A
999 Supreme Industrial Road
Anderson
South Carolina
12345
UNITED STATES

Payment Terms
1 Day Net

Delivery Date
3/31/19

For orders placed by AMCE US&A or AMCE US&B, the AMCE US&A / AMCE US&B Standard Conditions of Contract for the Purchase and / or Hire of Goods and Services (US&A/US&B) apply to this order.
For orders placed by AMCE US&A USA Corporation, the AMCE US&A USA Corporation Standard Conditions of

Use case: information extraction from tables

A
Company
Making
Everything

ACME US&A

999 Supreme Industrial Road , Anderson , South
Carolina .
12345 .UNITED STATES
Tel: +1-123-234-3456 Fed ID 11-123456
Fax: Tax Exempt 123345-6789
VAT Reg. No.:

Seller
We-Supply
NC Suppliers Corporation
P.O. Box 123456
Atlanta
GA 12345-98765
UNITED STATES
Phone: 111-222-3334
Fax:

Ship Via
Road

Terms Of Delivery
Free on board

Order

Our Order Date	Purchase Order
1/2/20	A12345
Revision	
1	

Supplier No
WESUPPLY

Delivery Address
ACME US&A
999 Supreme Industrial Road
Anderson
South Carolina

12345
UNITED STATES

Payment Terms
1 Day Net

Delivery Date
3/31/19

For orders placed by AMCE US&A or AMCE US&B, the AMCE US&A / AMCE US&B Standard Conditions of Contract for the Purchase and / or Hire of Goods and Services (US&A/US&B) apply to this order.
For orders placed by AMCE US&A USA Corporation, the AMCE US&A USA Corporation Standard Conditions of

Nodes:

- Each word is a node in the graph

Edges:

- Neighboring words are connected



Use case: information extraction from tables

ACME US&A
999 Supreme Industrial Road , Anderson , South Carolina .
12345 UNITED STATES
Tel: +1-123-234-3456 Fed ID 11-123456
Fax: Tax Exempt 123345-6789
VAT Reg. No.:

Order

Our Order Date	Purchase Order
1/2/20	A12345
Revision	
1	

Supplier No
WESUPPLY

Delivery Address
ACME US&A
999 Supreme Industrial Road
Anderson
South Carolina
12345
UNITED STATES

Payment Terms
1 Day Net

Delivery Date
3/31/19

Seller
We-Supply
NC Suppliers Corporation
P.O. Box 123456
Atlanta
GA 12345-98765
UNITED STATES
Phone: 111-222-3334
Fax:

Ship Via
Road

Terms Of Delivery
Free on board

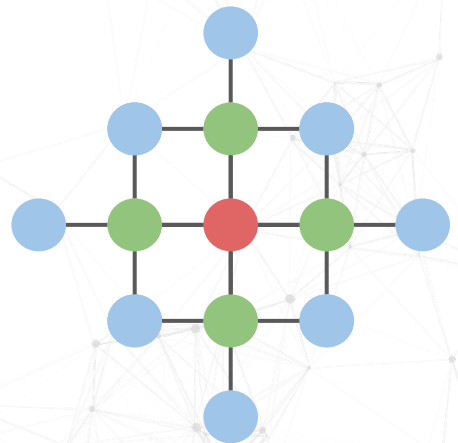
For orders placed by AMCE US&A or AMCE US&B, the AMCE US&A / AMCE US&B Standard Conditions of Contract for the Purchase and / or Hire of Goods and Services (US&A/US&B) apply to this order.
For orders placed by AMCE US&A USA Corporation, the AMCE US&A USA Corporation Standard Conditions of

Nodes:

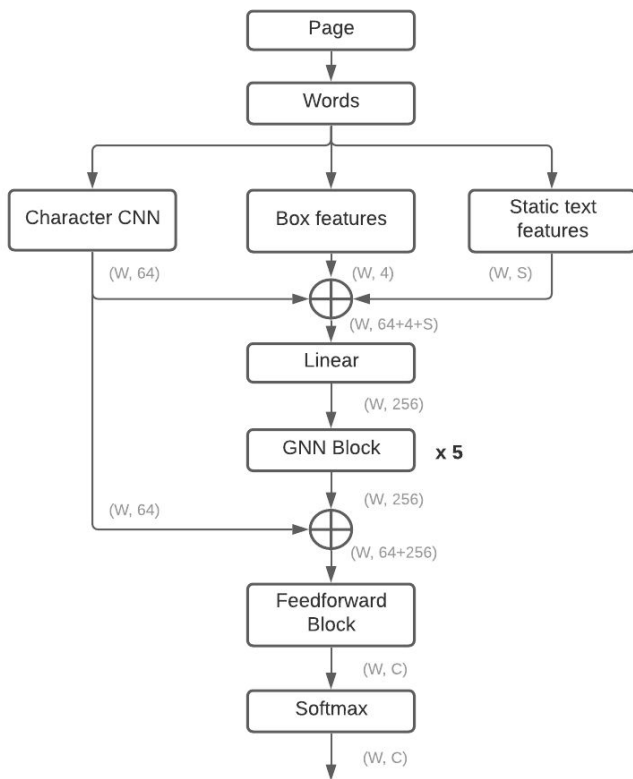
- Each word is a node in the graph

Edges:

- Neighboring words are connected



Use case: a possible model architecture



Legend

W: Number of words

S: Number of static text features

C: Number of classes

⊕: Concatenate tensors

Use case: results

ACME US&A
999 Supreme Industrial Road, Anderson, South Carolina, 12345 UNITED STATES
Tel: +1-23-234-3456 Fed ID 11-123456
Fax: Tax Exempt 123345-6789
VAT Reg. No.:

Purchase Order 012345-00
Vendor 000234

To :
We-Supply
DEPT AT 123456
ATLANTA GA 12345-45678
United States

Ship to :
ACME US&A
999 Supreme Industrial Road
ANDERSON, SC
United States

Phone () - Fax () -

PO Date	Ship Via	FOB	Planner	Confirming to	Terms
02/15/2019	Your Truck	OUR PLANT	DGJ		0.5% 10, NET 30
Item	Facility / Part / Rev / Description / Details	Vendor Quantity	Promised Delivery	Vendor Unit Cost	Extended Cost
1	Vendor U/M : CWT Vendor Desc : GALV, 12GA X 48-7/16" X 127-3/8" Default 12GA48716GR50 Rev 000 GALV, 12GA X 48-7/16" X 127-3/8" U/M EA Order Quantity: 11.0000	123.4567	02/18/2019	11.1111	
	Purchasing Category : INV **SHEET TOLERANCE** WIDTH +0/-0.15 LENGTH +0/-0.0625			222.2222	1,222.33
2	Vendor U/M : CWT Vendor Desc : GALV, 12GA X 48-7/16" X 138-3/4" SS Default 12GA48716138GR50 Rev 000 GALV, 12GA X 48-7/16" X 138-3/4" SS U/M EA Order Quantity: 12.0000	12.3456	02/18/2019	11.2222	
	Purchasing Category : INV **SHEET TOLERANCE** WIDTH +0/-0.15 LENGTH +0/-0.0625			333.4444	1,222.30
3		12.1111	02/15/2019	22.4444	



ACME US&A
999 Supreme Industrial Road, Anderson, South Carolina, 12345 UNITED STATES
Tel: +1-23-234-3456 Fed ID 11-123456
Fax: Tax Exempt 123345-6789
VAT Reg. No.:

Purchase Order 012345-00
Vendor 000234

To :
We-Supply
DEPT AT 123456
ATLANTA GA 12345-45678
United States

Ship to :
ACME US&A
999 Supreme Industrial Road
ANDERSON, SC
United States

Phone () - Fax () -

PO Date	Ship Via	FOB	Planner	Confirming to	Terms
02/15/2019	Your Truck	OUR PLANT	DGJ		0.5% 10, NET 30
Item	Facility / Part / Rev / Description / Details	Vendor Quantity	Promised Delivery	Vendor Unit Cost	Extended Cost
1	Vendor U/M : CWT Vendor Desc : GALV, 12GA X 48-7/16" X 127-3/8" Default 12GA48716GR50 Rev 000 GALV, 12GA X 48-7/16" X 127-3/8" U/M EA Order Quantity: 11.0000	123.4567	02/18/2019	11.1111	
	Purchasing Category : INV **SHEET TOLERANCE** WIDTH +0/-0.15 LENGTH +0/-0.0625			222.2222	1,222.33
2	Vendor U/M : CWT Vendor Desc : GALV, 12GA X 48-7/16" X 138-3/4" SS Default 12GA48716138GR50 Rev 000 GALV, 12GA X 48-7/16" X 138-3/4" SS U/M EA Order Quantity: 12.0000	12.3456	02/18/2019	11.2222	
	Purchasing Category : INV **SHEET TOLERANCE** WIDTH +0/-0.15 LENGTH +0/-0.0625			333.4444	1,222.30
3		12.1111	02/15/2019	22.4444	



About 1000 training documents



Macro F1 score: 88%

Implementations, literature

- [PyTorch Geometric](#)
 - “[GNN Cheatsheet](#)” in the docs contains an interesting list of papers
- [Spektral](#) for TensorFlow
- [Deep Graph Library](#) (framework agnostic)
- https://en.wikipedia.org/wiki/Graph_neural_network exists since 5 July 2021