

# 辺編集による文書レベルの関係グラフ構築

牧野晃平 三輪誠 佐々木裕

豊田工業大学

{sd19446, makoto-miwa, yutaka.sasaki}@toyota-ti.ac.jp

## 1 はじめに

文献中に記述された情報を抽出して利用する研究が盛んに行われている [1, 2, 3]. 抽出対象の情報は、表形式や用語ペアと関係の集合などのように、目的に応じて様々な形式をとる. これらの情報の抽出にはその形式に適した手法が必要だが、それぞれに特化した手法を作成するのは非効率でコストが高い.

情報を表現する形式の一つとして、グラフ構造がある. グラフ構造は、節点と、節点同士を関連性に基づいて結びつける辺の集合で表現できる形式で、様々な情報を包括的に表現可能な形式の一つである. 例えば、表形式であれば、見出しと各セルを節点として、その間を辺で接続して表現できる. このグラフ構造を利用して、文献中から抽出する情報をグラフ構造として扱うことで、様々な情報抽出タスクに利用可能な抽出器が実現できると考えられる.

そこで本研究では、上記の抽出器作成の第一歩として、用語が与えられた上での文を超えた文書レベルでの関係分類を対象に、節点を与えられた上でのグラフの構築モデルの実現を目指す. 対象の用語間の関係を、用語を節点・関係を辺としたグラフと捉えることで対象問題をグラフ構造として扱う.

提案手法では、他のシステムで構築されたグラフをもとに、辺を徐々に編集する深層学習モデルを実現する. 既存のグラフの辺を徐々に編集することで、それまでに抽出されている情報をグラフとして利用し、文脈情報と周囲のグラフ情報を考慮したグラフを構築する.

## 2 関連研究

### 2.1 グラフ畳み込みネットワーク

Kipf らはグラフ構造の表現獲得のためにグラフ畳み込みネットワーク (Graph Convolutional Network; GCN) を提案した [4]. GCN では、各节点の特徴を

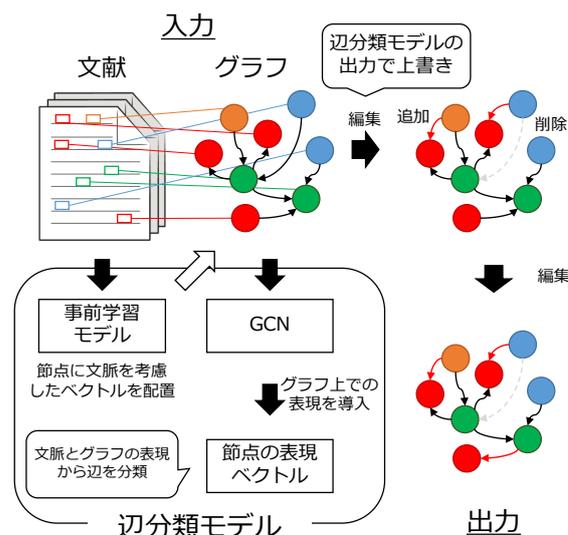


図1 提案手法の全体像

表すベクトルを用いて、節点の表現を周辺の節点の情報を用いて更新することで、グラフ構造における周辺の情報を考慮した各節点の表現を得る.

### 2.2 情報抽出

グラフを抽出する類似のタスクとして、関係抽出や時間関係抽出などがある. これらのタスクでは、入力された文中の用語間の関係性を抽出する. Christopoulou ら [1] は用語を節点・関係を辺とした文書単位のグラフを作成し、その辺を分類することで、文間の関係を含む文書単位の関係抽出に取り組んだ. Nan ら [5] は文書単位で用語や単語の潜在的なグラフを構築することで、関係抽出を行った. Cheng ら [6] は事象それぞれに対して時刻にあわせて動的な表現を利用して時間関係抽出を行った.

## 3 グラフの辺編集モデル

本研究では、用語が与えられた上での文を超えた文書レベルでの関係分類をグラフの構築として定式化する. 具体的には、図1のように節点を与えられた上で既存のシステムで抽出したグラフの辺を編集

して、文脈上とグラフ上の両者の観点を考慮に入れたグラフを構築する。グラフの辺の編集は、文脈情報とグラフ構造を考慮したグラフの辺分類モデルによって節点ペア間の辺のクラスを分類し、その結果で上書きすることで実現する。

モデルは文脈情報を導入するために事前学習モデル Longformer [7] を利用して節点の表現を作成する。また、グラフ上で周辺の節点の情報を導入するために、節点の情報を接続している辺に沿って伝搬させる手法である GCN を利用する。

### 3.1 グラフの辺分類モデル

グラフの辺分類モデルとして、グラフの任意の節点ペア間の辺を分類するモデルを提案する。モデルは、入力として文書の情報  $\text{doc}$  と、節点が  $\mathcal{N}$ ・辺が  $\mathcal{E}$  のグラフ、編集対象の節点ペア  $(\mathcal{N}_i, \mathcal{N}_j)$  を受け取り、編集対象の節点ペア間の辺のクラス  $\mathcal{E}_{ij}$  を出力する。ここでは、文書  $\text{doc}$  と節点に対応する用語の情報  $\mathcal{N}$  から各節点の文脈を考慮した表現を得る  $\text{EncodeNode}(\text{doc}, \mathcal{N})$  と、得られた節点の表現  $\vec{\mathcal{N}}$  を用いて編集対象の辺の表現  $\vec{\mathcal{E}}$  を作成する  $\text{EncodeEdge}(\vec{\mathcal{N}}, \mathcal{E})$ 、そして節点ペア  $(\mathcal{N}_i, \mathcal{N}_j)$  間の辺の表現  $\vec{\mathcal{E}}_{ij}$  から編集後の辺  $\hat{\mathcal{E}}_{ij}$  を出力する  $\text{EdgeClassifier}(\vec{\mathcal{E}}_{ij})$  に分けて説明する。これらを合わせると以下の式で表せる。

$$\vec{\mathcal{N}} = \text{EncodeNode}(\text{doc}, \mathcal{N}) \quad (1)$$

$$\vec{\mathcal{E}} = \text{EncodeEdge}(\vec{\mathcal{N}}, \mathcal{E}) \quad (2)$$

$$\hat{\mathcal{E}}_{ij} = \text{EdgeClassifier}(\vec{\mathcal{E}}_{ij}) \quad (3)$$

各節点の表現を作成する  $\text{EncodeNode}$  は、事前学習した Longformer [7] から得られたサブワード単位の表現を用語レベルの表現にまとめるために、サブワードの表現の各次元で最大値をとるプーリング (Pool) をして、用語のラベルの表現  $\mathbf{v}^{lab}$  を結合して表現する。

$$\begin{aligned} \vec{\mathcal{N}} &= \text{EncodeNode}(\text{doc}, \mathcal{N}) \\ &= [\text{Pool}(\text{Longformer}(\text{doc})); \mathbf{v}^{lab}] \end{aligned} \quad (4)$$

$\text{EncodeEdge}(\vec{\mathcal{N}}, \mathcal{E})$  では、 $\text{EncodeNode}$  によって得られた節点の表現  $\vec{\mathcal{N}}$  から、辺の表現  $\vec{\mathcal{E}}$  を作成する。まず、節点の表現を、GCN によって辺に沿って伝搬させることで、グラフ上での周辺の節点を考慮した節点の表現  $\vec{\mathcal{N}}^G$  を計算する。

$$\vec{\mathcal{N}}^G = \text{GCN}(\vec{\mathcal{N}}, \mathcal{E}) \quad (5)$$

---

#### Algorithm 1: グラフの編集手順

---

$\text{Distance}(\mathcal{N}, d_1, d_2)$ : 距離が  $d_1$  以上  $d_2$  未満の節点ペア集合  $\mathcal{P}$  を返す関数

**Input:**  $\text{doc}$ : 文書,  $\mathcal{N}$ : 節点の集合,  $\mathcal{E}$ : 辺集合,  $d_{max}$ : 距離の最大値

**Output:**  $\mathcal{E}$ : 辺集合

$\vec{\mathcal{N}} \leftarrow \text{EncodeNode}(\text{doc}, \mathcal{N})$

**while**  $d$  in range( $\max(|\mathcal{N}|, d_{max})$ ) **do**

$\vec{\mathcal{E}} \leftarrow \text{EncodeEdge}(\vec{\mathcal{N}}, \mathcal{E})$

**if**  $d = d_{max}$  **then**

$\mathcal{P} \leftarrow \text{Distance}(\mathcal{N}, d_{max}, \infty)$

**else**

$\mathcal{P} \leftarrow \text{Distance}(\mathcal{N}, d, d + 1)$

**end if**

**while**  $(i, j)$  in  $\mathcal{P}$  **do**

$\vec{\mathcal{E}}_{ij} \leftarrow \text{EdgeClassifier}(\vec{\mathcal{E}}_{ij})$

**end while**

**end while**

---

各辺の表現  $\vec{\mathcal{E}}$  の各要素として、各節点の組み合わせ  $(\mathcal{N}_i, \mathcal{N}_j)$  について、その間の辺の表現ベクトル  $\vec{\mathcal{E}}_{ij}$  を計算する。この表現は、グラフから得られた表現に、分類に有効となるような人手の特徴ベクトルとして、隣接節点を表す表現や節点の始点と終点が本文中でどちらが先に出現するかを表す表現  $\mathbf{b}_{ij}$  を結合して作成する。辺の始点と終点となる節点に対する全結合層をそれぞれ、 $\text{FC}^{head}$  と  $\text{FC}^{tail}$  とすると、以下のように表現できる。

$$\vec{\mathcal{E}} = \text{EncodeEdge}(\vec{\mathcal{N}}, \mathcal{E}) \quad (6)$$

$$\vec{\mathcal{E}}_{ij} = [\text{FC}^{head}(\vec{\mathcal{N}}_i^G) \otimes \text{FC}^{tail}(\vec{\mathcal{N}}_j^G); \mathbf{b}_{ij}] \quad (7)$$

ただし、 $\otimes$  は要素積を示す。

$\text{EdgeClassifier}(\vec{\mathcal{E}}_{ij})$  では得られた節点間の表現から、節点  $(\mathcal{N}_i, \mathcal{N}_j)$  間の辺  $\mathcal{E}_{ij}$  を分類する。出力の全結合層  $\text{FC}^{out}$  を用いて、それぞれのクラスの確率  $\hat{\mathbf{p}}_{ij}$  を計算し、最大の確率のクラスを選択して、辺の予測値  $\hat{\mathcal{E}}_{ij}$  を出力する。

$$\hat{\mathcal{E}}_{ij} = \text{EdgeClassifier}(\vec{\mathcal{E}}_{ij}) = \arg \max \hat{\mathbf{p}}_{ij} \quad (8)$$

$$\hat{\mathbf{p}}_{ij} = \text{Softmax}(\text{FC}^{out}(\vec{\mathcal{E}}_{ij})) \quad (9)$$

### 3.2 グラフの編集

グラフの構築は、入力されたグラフの節点ペアの全組み合わせに対して、3.1 節の辺分類モデルによる分類結果によって徐々に上書きして行くことで実

現する。編集の順序は、近い関係を先に抽出し、遠い問題は後に抽出する近傍優先戦略 [8, 9] に基づき、節点の元の用語同士が文書内で近いものから編集を進め、順次、離れた節点同士の編集を行う。用語の距離は出現順序に基づいて計算する。例えば文書中で  $m$  番目に出現した用語と  $m+3$  番目に出現した用語同士の距離は 3 とする。

具体的な編集方法をアルゴリズム 1 に示した。距離が近い節点ペアから順番に、同一距離のペアごとにまとめて、辺分類器による辺のクラス分類を行い、その出力で上書きをして編集を進める。このように編集を進めることで、すべての辺を独立に編集する場合では  $|N|^2$  回の編集が必要であるのに比べて、 $|N|-1$  回の編集で行うことができ、さらに上限  $d_{max}$  を設けることで、編集は  $d_{max}$  回で完了する。

### 3.3 学習

辺分類モデルの学習は、正しいグラフが抽出できるように、対数尤度を最大化する。正解の辺を  $\mathcal{E}_{ij}^{gold}$ 、辺の全クラスを  $\mathcal{C}$  とすると、損失  $\mathcal{L}$  は以下のように定義できる。

$$\mathcal{L} = - \sum_{i=1}^{|N|} \sum_{j=1}^{|N|} \sum_{c \in \mathcal{C}} \mathbb{1}[c = \mathcal{E}_{ij}^{gold}] \cdot w_c \log \hat{p}_{ij}[c] \quad (10)$$

ただし、 $\mathbb{1}[\cdot]$  は、括弧内の条件を満たしたときに 1、満たさないときには 0 を返す関数で、 $w_c$  はクラスごとの重みである。クラスごとの重み  $w_c$  は、それぞれのクラスの損失の重みを示す。

## 4 実験と考察

### 4.1 実験設定

本実験では無機材料文献からの合成プロセス抽出のためのコーパス合成手順コーパス [10] を利用して実験を行う。合成手順コーパスでタグ付けされている合成プロセスは、材料や操作、条件などを表す用語と、操作の進行や材料の投入、条件付けといったような関係がタグ付けされているコーパスで、各文献に対して一つのグラフを定義することができる。訓練データ・開発データ・評価データはそれぞれ  $200 \cdot 15 \cdot 15$  で、公開されているコーパスと同一の分割で実験を行う。

評価指標は、全体の F 値であるマイクロ F 値と、各クラスの F 値を平均したマクロ F 値で評価する。評価は開発データに対して最高のマイクロ F 値が得

表 1 探索するパラメタ

ハイパパラメタ名	探索対象のリスト
GCN の層の数	0, 1, 2, 3
次元数	64, 128, 256
$d_{max}$	1, 5, 10, 15

表 2 グラフの編集モデルとルールベース抽出器の比較

モデル	マイクロ F	マクロ F
ルール+編集 (提案手法)	0.802	<b>0.788</b>
ルールベース	<b>0.807</b>	0.689
編集のみ	0.795	0.720
GCN 無し	0.775	0.728

られた点での評価データに対する評価値で行う。

作成したルールの抽出結果を編集した提案手法 (ルール+編集) に対する比較手法として、ルールによる抽出 (ルールベース)、辺の全く接続されていないグラフから編集によって辺を接続していくもの (編集のみ)、そして (7) 式において  $\bar{N}^G$  の代わりに  $\bar{N}$  を用いた GCN を利用せずにすべての辺を同時に決定したもの (GCN 無し) の三つを用意した。ルールの詳細は付録 A に示した。ルールは節点ペアのクラスの組み合わせごとに設定し、最も近いもの同士で接続したり、出現した順番に接続したり、辞書マッチを用いたりした。

表 1 に示したハイパパラメタについては、パラメタの組み合わせを全探索し、最適なパラメタとして下線が引かれたパラメタを選択した。全探索はどのモデルも十分収束する 130 エポックで学習を行い、モデルの学習が終了したとき、開発データに対して最も高いマイクロ F 値を記録したエポックのときのモデルをその試行のモデルとした。パラメタは開発データに対して最大のマイクロ F 値を記録したパラメタを選択した。

事前実験によって損失における  $w_c$  を  $w_c = 1$  として学習した場合では、予測する辺の数のうちのほとんどが負例であることから、モデルが負例しか予測しなくなってしまったため、正例のクラスの損失が大きくなるようにした。学習初期には正例の損失が大きくなるように係数を掛けておき、学習が進むにつれ負例と同一の重み付けになるように係数をスケジューリングした。

### 4.2 実験結果

得られた結果を表 2、ルールベースとルール+編集に対するクラスごとの抽出性能は付録 B に示し



## 参考文献

- [1] Fenia Christopoulou, Makoto Miwa, and Sophia Ananiadou. Connecting the dots: Document-level neural relation extraction with edge-oriented graphs. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pp. 4924–4935, 2019.
- [2] Makoto Miwa and Mohit Bansal. End-to-end relation extraction using LSTMs on sequences and tree structures. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1105–1116, Berlin, Germany, August 2016. Association for Computational Linguistics.
- [3] Patrick Verga, Emma Strubell, and Andrew McCallum. Simultaneously self-attending to all mentions for full-abstract biological relation extraction. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 872–884, New Orleans, Louisiana, June 2018. Association for Computational Linguistics.
- [4] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [5] Guoshun Nan, Zhijiang Guo, Ivan Sekulic, and Wei Lu. Reasoning with latent structure refinement for document-level relation extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1546–1557, Online, July 2020. Association for Computational Linguistics.
- [6] Fei Cheng, Masayuki Asahara, Ichiro Kobayashi, and Sadao Kurohashi. Dynamically updating event representations for temporal relation classification with multi-category learning. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 1352–1357, Online, November 2020. Association for Computational Linguistics.
- [7] Iz Beltagy, Matthew E. Peters, and Arman Cohen. Longformer: The long-document transformer. *arXiv:2004.05150*, 2020.
- [8] Makoto Miwa and Yutaka Sasaki. Modeling joint entity and relation extraction with table representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1858–1869, Doha, Qatar, October 2014. Association for Computational Linguistics.
- [9] Shuai Ma, Gang Wang, Yansong Feng, and Jinpeng Huai. Easy first relation extraction with information redundancy. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3851–3861, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [10] Sheshera Mysore, Zachary Jensen, Edward Kim, Kevin Huang, Haw-Shiuan Chang, Emma Strubell, Jeffrey Flanigan, Andrew McCallum, and Elsa Olivetti. The materials science procedural text corpus: Annotating materials synthesis procedures with shallow semantic structures. In *Proceedings of the 13th Linguistic Annotation Workshop*, pp. 56–64, Florence, Italy, August 2019. Association for Computational Linguistics.
- [11] S.B. Zhang, Y.P. Sun, B.C. Zhao, X.B. Zhu, and W.H. Song. Influence of ni doping on the properties of perovskite molybdates  $\text{SrMo}_{1-x}\text{Ni}_x\text{O}_3$  ( $0.02 \leq x \leq 0.08$ ). *Solid State Communications*, Vol. 142, No. 12, pp. 671 – 675, 2007.

## A ルールベースの抽出器

合成手順コーパスの関係を抽出するためのルールベースの抽出器では、記述されたルールに従い、ルールに適合するものに対して辺を接続する。辺の種類は OPERATION-OPERATION・OPERATION-MATERIAL・その他の辺の3種類に分類してルールを定義する。ルールは節点ペアの用語ラベルごとに出現順序と距離に合わせて定義する。距離は用語ペア間に存在する語の数で定義し、文書をスペース区切りにすることで単語分割した。

### A.1 OPERATION-OPERATION

NEXT\_OPERATION : OPERATION-OPERATION の辺ラベルは、NEXT\_OPERATION のみの一種類で、OPERATION の用語が出現した順番に、前から後ろへ NEXT\_OPERATION の辺をつける。

### A.2 OPERATION-MATERIAL

SOLVENT\_MATERIAL・ATMOSPHERIC\_MATERIAL・PARTICIPANT\_MATERIAL : これらのクラスについては、クラス毎に辞書を用意し、マッチした辞書に割り当てられたクラスの辺を MATERIAL と最近傍の OPERATION に対する辺を接続する。

RECIPE\_PRECURSOR・RECIPE\_TARGET : SOLVENT\_MATERIAL, ATMOSPHERIC\_MATERIAL, PARTICIPANT\_MATERIAL の辞書に該当しない MATERIAL をすべて、実際に訓練データに存在する辺の数がより多い RECIPE\_PRECURSOR のクラスの辺として最近傍の OPERATION に対して接続する。

### A.3 その他の関係

PROPERTY\_OF : PROPERTY-UNIT を始点とする場合では、PROPERTY-UNIT から文内で最も近傍に存在する MATERIAL に対して PROPERTY\_OF の辺を割り当てる。PROPERTY-MISC を始点とする場合では、文中で PROPERTY-MISC から最も近傍に存在する MATERIAL もしくは NONRECIPE-MATERIAL に対して接続する。

CONDITION\_OF : すべての CONDITION-UNIT と CONDITION-MISC から最近傍の OPERATION に対して、CONDITION\_OF の関係を割り当てる。

NUMBER\_OF : NUMBER から文内の後方で PROPERTY-UNIT・CONDITION-UNIT・APPARATUS-UNIT に属する用語の内、最も近傍に記述されている用語に対して辺を接続する。

AMOUNT\_OF : AMOUNT-UNIT と AMOUNT-UNIT から、MATERIAL と NONRECIPE-MATERIAL のうち、文内で最近傍のものに対して辺を接続する。

DESCRIPTOR\_OF : MATERIAL-DESCRIPTOR から、文内で最も近傍に存在する MATERIAL もしくは NONRECIPE-MATERIAL に対して、DESCRIPTOR\_OF の辺を接続する。また、APPARATUS-DESCRIPTOR からの関係では、SYNTHESIS-APPARATUS に対してのみ辺を接続する。

APPARATUS\_OF : (SYNTHESIS-APPARATUS|CHARACTERIZATION-APPARATUS)-OPERATION の辺では、SYNTHESIS-APPARATUS 及び CHARACTERIZATION-APPARATUS から、前方で最も近傍に存在する OPERATION に対して辺を接続する。前方に存在しない場合は、後方で最も近傍に存在する OPERATION に対して辺を接続する。

TYPE\_OF : まず、PROPERTY-TYPE-PROPERTY-UNIT と、APPARATUS-PROPERTY-TYPE-APPARATUS-UNIT については、

それぞれ文中で最も近傍に存在する単位に対して関係を割り当てる。CONDITION-TYPE-CONDITION-UNIT については、CONDITION-TYPE から文中の前方で最も近傍に存在する CONDITION-UNIT に関係を割り当てる。

BRAND\_OF : BRAND から前方で MATERIAL・NONRECIPE-MATERIAL・SYNTHESIS-APPARATUS・CHARACTERIZATION-APPARATUS のクラスの用語の内、文内で最も近傍に存在するものに対して辺を接続する。

APPARATUS\_ATTR\_OF : APPARATUS-UNIT から最近傍の SYNTHESIS-APPARATUS もしくは CHARACTERIZATION-APPARATUS に対して辺を接続する。

COREF\_OF : ルールを記述するのは困難であったため、ルールベースでは対象としない。

## B クラスごとの抽出性能

ルールベースとルール+編集のクラスごとに算出した評価値を、評価データに対して算出したものをそれぞれ表3と表4に示した。

表3 ルールベースの抽出器におけるクラス毎の抽出結果

辺クラス	Precision	Recall	F 値
NEXT_OPERATION	0.990	0.881	0.932
RECIPE_PRECURSOR	0.730	0.414	0.528
RECIPE_TARGET	0.000	0.000	0.000
PARTICIPANT_MATERIAL	0.419	0.800	0.550
SOLVENT_MATERIAL	0.697	0.418	0.522
ATMOSPHERIC_MATERIAL	1.000	0.378	0.549
PROPERTY_OF	0.905	1.000	0.950
CONDITION_OF	0.963	0.981	0.972
NUMBER_OF	0.943	0.961	0.952
AMOUNT_OF	0.744	0.865	0.800
DESCRIPTOR_OF	0.931	0.979	0.955
TYPE_OF	0.769	1.000	0.870
BRAND_OF	0.561	0.920	0.697
APPARATUS_OF	0.972	0.854	0.909
APPARATUS_ATTR_OF	0.909	0.769	0.833
COREF_OF	0.000	0.000	0.000

表4 ルール+編集におけるクラスごとの抽出性能

辺クラス	Precision	Recall	F 値
NEXT_OPERATION	0.822	0.738	0.778
RECIPE_PRECURSOR	0.517	0.605	0.558
RECIPE_TARGET	0.818	0.692	0.750
PARTICIPANT_MATERIAL	0.524	0.619	0.568
SOLVENT_MATERIAL	0.788	0.473	0.591
ATMOSPHERIC_MATERIAL	0.929	0.867	0.897
PROPERTY_OF	0.905	0.864	0.884
CONDITION_OF	0.925	0.892	0.908
NUMBER_OF	0.990	0.976	0.983
AMOUNT_OF	0.769	0.869	0.816
DESCRIPTOR_OF	0.941	0.914	0.928
TYPE_OF	0.923	0.923	0.923
BRAND_OF	0.634	0.897	0.743
APPARATUS_OF	0.750	0.675	0.711
APPARATUS_ATTR_OF	0.909	0.667	0.769
COREF_OF	0.714	0.909	0.800