# Knowledge Graph Embeddings for Downstream Machine Learning Tasks with RDF2Vec

Terencio AGOZZINO

`terencio.agozzino@std.heh.be`

MA2 IT Engineering

*Advisor:*
Prof. Dr. Femke ONGENAE

*Supervisors:*
Dr. Ir. Gilles VANDEWIELE
Dr. Ir. Samuel CREMER
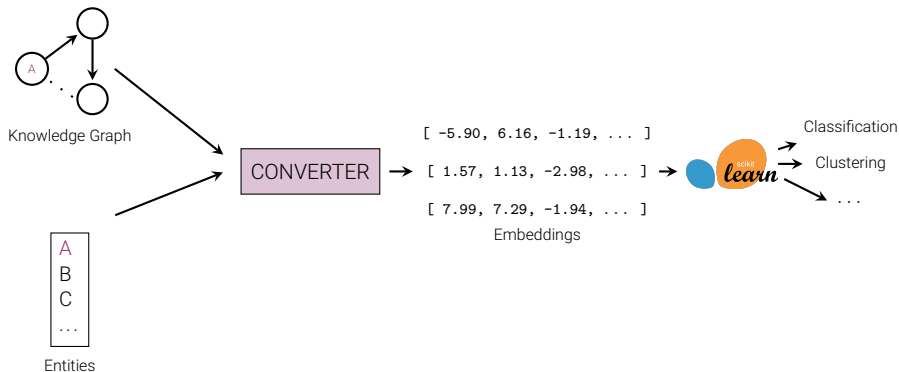Ir. Bram STEENWINCKEL

# IDLab

# IDLab[1]

- Influential research group of imec.
- Focuses on Internet technologies and data science.
- Makes contributions to future global standards (e.g., W3C and 5G).
- Located at the University of Ghent and the University of Antwerp.
- Counts more than 500 collaborators.

---

[1]Internet Technology & Data Science Lab.

# Internship

- Associated with the Master's Thesis over a period of 16 weeks.
- Took place within the Knowledge Management team of the IDLab's Discover research center cluster which optimizes intelligent agents and their interactions.
- The Knowledge Management team created more than twenty research projects.
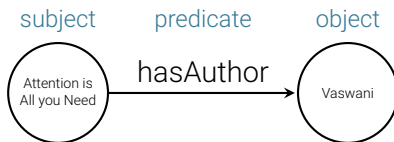
# Motivation

# Motivation



Knowledge Graph

Entities

CONVERTER

[ -5.90, 6.16, -1.19, ... ]
[ 1.57, 1.13, -2.98, ... ]
[ 7.99, 7.29, -1.94, ... ]
Embeddings

Classification
Clustering
...

# Knowledge Graph

Figure: Triple.

Example of URI[2]: http://dl-learner.org/carcinogenesis #d187

$URL^{3}$     $URN^{4}$

---

[2]Uniform Resource Identifier.
[3]Uniform Resource Locator.
[4]Uniform Resource Name.

# Knowledge Graph

subject     predicate     object

Figure: Triple.

Example of URI[2]: $\underbrace{\texttt{http://dl-learner.org/carcinogenesis}}_{URL^3}\underbrace{\texttt{\#d187}}_{URN^4}$

---

[2]Uniform Resource Identifier.
[3]Uniform Resource Locator.
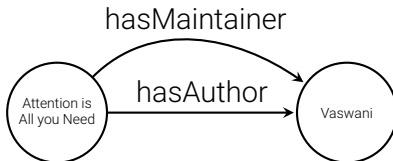[4]Uniform Resource Name.
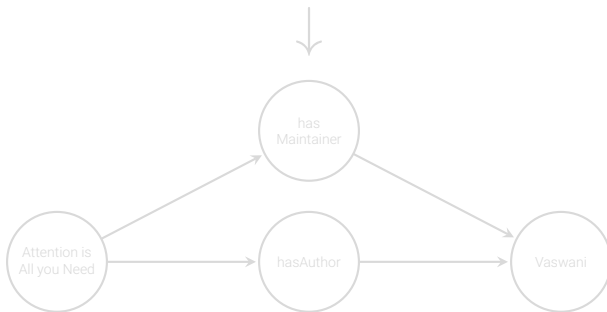
# Knowledge Graph

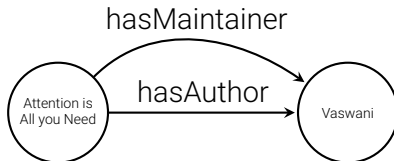Figure: Knowledge Graph (KG).



Figure: Oriented Graph.

# Knowledge Graph

Figure: Knowledge Graph (KG).
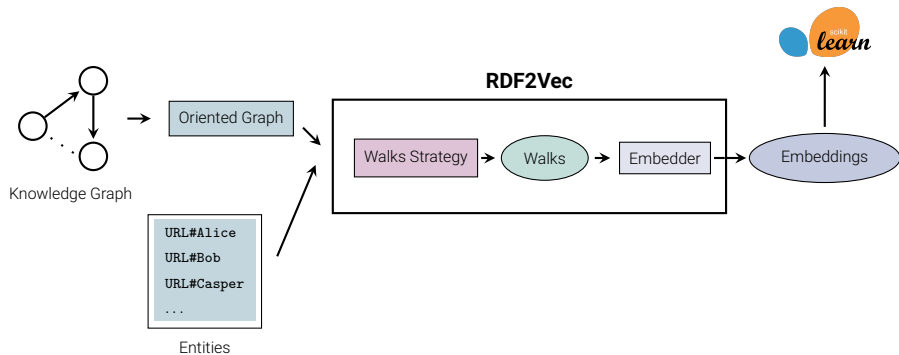


Figure: Oriented Graph.

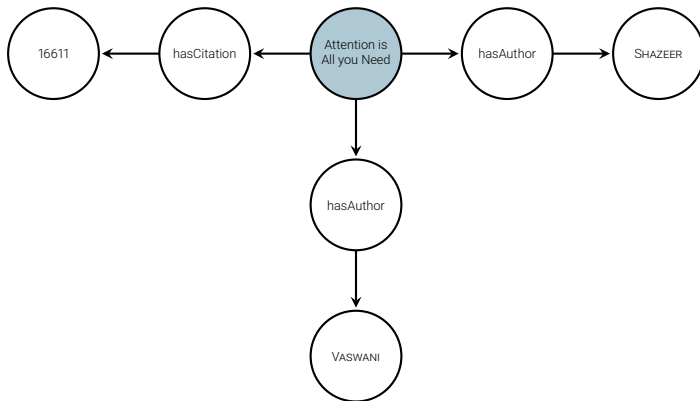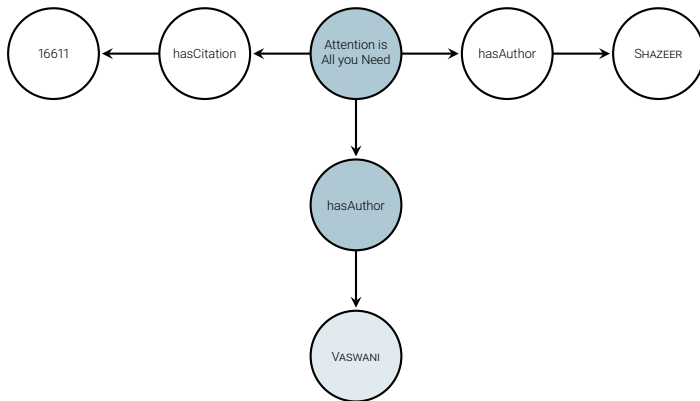# RDF2Vec[5]

---

# RDF2Vec

# RDF2Vec

Figure: Walk Extraction for an Oriented Graph (Part I).

Walk Extraction: Attention is All you Need

# RDF2Vec

Figure: Walk Extraction for an Oriented Graph (Part II).

Walk Extraction: Attention is All you Need → hasAuthor

# RDF2Vec
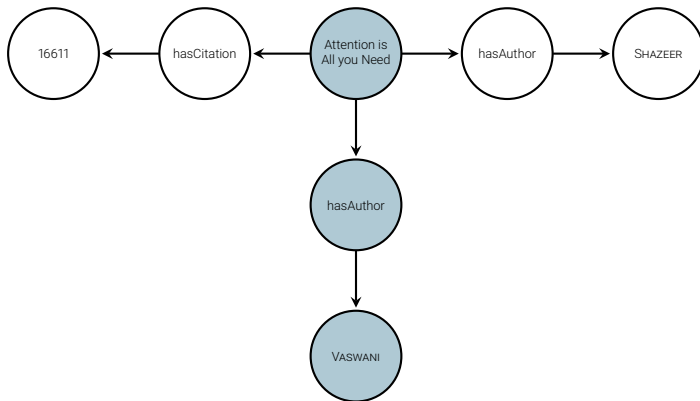
Figure: Walk Extraction for an Oriented Graph (Part III).

Walk Extraction: Attention is All you Need → hasAuthor → VASWANI

# Internship

# Scope of the Internship

Consists in improving `pyRDF2Vec`:

1. Add better support for different KG sizes.
2. Add first support of *literals*.
3. Avoid re-training a model when new nodes[6] are added in the KG.
4. Improve the documentation of the library.
5. Add Continuous Integration / Delivery / Testing.

Where the goals are *ordered according to importance*.
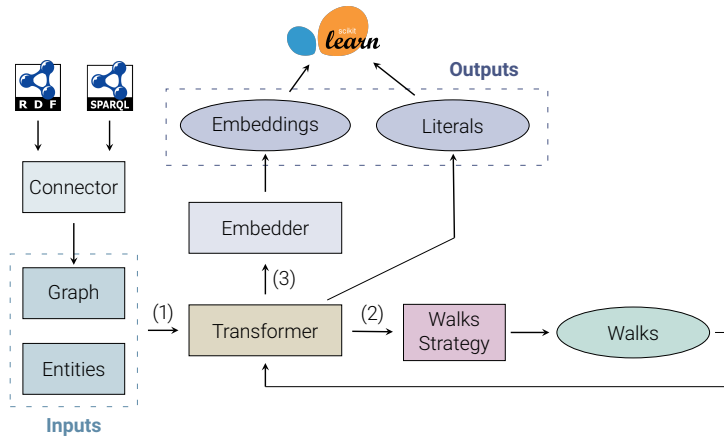
---

[6]Also called *vertices*.

# Architecture



Figure: Workflow of `pyRDF2Vec`.

# Optimization Mechanisms

Consists of reducing the walk extraction time by using:

- a cache memory;
- multiprocessing;
- cache pre-filling;
- an appropriate data structure for the nodes;
- a connection pool;
- …

# Optimization Mechanism Results

# Benchmarks Setup

| Data Set | Triples | Entities | Relations | Size |
|----------|---------|----------|-----------|------|
| `MUTAG` | 74,567 | 22,534 | 24 | small |
| `AM` | 5,700,371 | 933,471 | 100 | medium |
| `DBP:Cities` | 651,580,976 | 4,233,000 | 7,992 | large |

Table: Properties of the Data Sets Used for the Benchmarks.

Performed on the IDLab's servers with 4 CPUs and 64 GB RAM.

# Benchmark Results

| Data Set | Speedup | | |
|:---:|:---:|:---:|:---:|
| | **Minimum** | **Maximum** | **Average** |
| MUTAG | 6.20 | 13.28 | 9.74 |
| AM | 1.44 | 8.94 | 5.19 |
| DBP:Cities | 2.22 | 3.23 | 2.72 |

Table: Speedup Results with Optimization Mechanisms.

- caching and pre-filling benefit small KGs like MUTAG the most;
- multiprocessing is more interesting for large KGs like DBP:Cities.
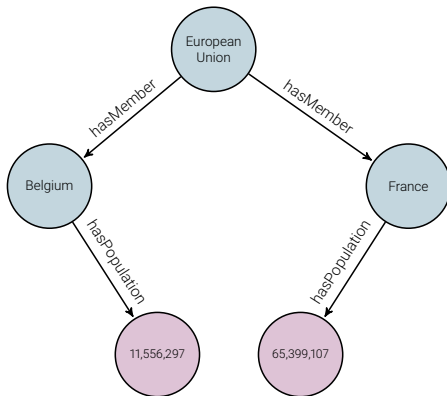
# Literals

# Literals

Figure: Oriented Graph with 5 Nodes and 4 Edges.

# Literals

Literals may:

- be specific to some entities;
- be too expensive to extract;
- not be interesting.

→ Let the user specify the literals to extract

For example: [

```
                predicate 1        predicate 2
             ⎴⎴⎴⎴⎴⎴⎴⎴⎴         ⎴⎴⎴⎴⎴⎴⎴⎴⎴⎴⎴
        [ hasCapital → hasLatitude ],    # 50.85; 48.85
        [ hasCapital → hasLongitude ],   # 4.35; 2.34
        [ hasPopulation ],               # 11,55; 65,39
        ...

        ]
```

# Literals

Literals may:

- be specific to some entities;
- be too expensive to extract;
- not be interesting.

→ Let the user specify the literals to extract

For example: [

```
            predicate 1      predicate 2
       [ hasCapital → hasLatitude ],   # 50.85; 48.85
       [ hasCapital → hasLongitude ],  # 4.35; 2.34
       [ hasPopulation ],              # 11,55; 65,39
       ...

     ]
```

# Online Learning

# Online Learning

Figure: Oriented Graph With 7 Nodes and 6 Edges.

→ Consists in updating the vocabulary of walks of an initial model.

# Online Learning

Figure: Oriented Graph With 7 Nodes and 6 Edges.

→ Consists in updating the vocabulary of walks of an initial model.

# Demo Time! 🥁

On the menu:

- The `pyRDF2Vec` GitHub repository.
- An example of using literals.
- ~~An example of using online learning~~ (available on GitHub).

# A Trip to Sesame Street: Evaluation of BERT and Other Recent Embedding Techniques Within RDF2Vec

# RDF2Vec

# RDF2Vec

# RDF2Vec

Figure: Oriented Graph.

# RDF2Vec

Figure: Walk Extraction for an Oriented Graph (Part I).

Walk Extraction: Attention is All you Need

# RDF2Vec

Figure: Walk Extraction for an Oriented Graph (Part II).

Walk Extraction: Attention is All you Need → hasAuthor

# RDF2Vec

Figure: Walk Extraction for an Oriented Graph (Part III).

Walk Extraction: Attention is All you Need → hasAuthor → VASWANI

# RDF2Vec

# RDF2Vec

# Scope of the Master's Thesis

Mainly consists in evaluating BERT[7] with Word2Vec and implementing other strategies:

1. Supported BERT and FastText for comparison purposes.
2. Evaluated the impact of BERT with Word2Vec and FastText.
3. `SplitWalker` as a new walking strategy in RDF2Vec.
4. `WideSampler` as a new sampling strategy in RDF2Vec.

`BONUS`: has unintentionally improved Word2Vec and `pyRDF2Vec`.

---

[7]Only the traditional BERT model is concerned by this scope.

# BERT

# BERT

Figure: Summary Timeline.

---

[8]Recurrent Neural Networks.

# BERT

Main characteristics:

- Allows one or two input sentences.
- Uses a WordPiece tokenization.

Table: Example of Tokenization With BERT.

|  | The | machine | loves | embeddings. |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|
| [CLS] | the | machine | loves | em | ##bed | ##ding | ##s | . | [SEP] |
| 101 | 1996 | 3698 | 7459 | 7861 | 8270 | 4667 | 2015 | 1012 | 102 |

- Uses two pre-training phases:
  1. Masked Language Model (MLM);
  2. Next Sentence Prediction (NSP).
- Contextualizes embeddings using bidirectional representations.

# BERT Implementation[9] With KGs

- Build the vocabulary.
- Fit the BERT model in two ways:
  1. node tokenization: ensuring not to split the nodes;
  2. training: only done with MLM using a data collator.
- Get entity embeddings.

[9] https://github.com/IBCNServices/pyRDF2Vec/blob/feature/bert/pyrdf2vec/embedders/bert.py

# Improve the Model's Accuracy of Word2Vec

# Improve the Model's Accuracy of Word2Vec

Table: Context Words Determination for a Window Size of 2.

| Input Text | Target Word | Context Words |
|---|---|---|
| I will always remember her | i | will<br>always |
| I will always remember her | will | i<br>always<br>remember |
| I will always remember her | always | i<br>will<br>remember<br>her |
| I will always remember her | remember | will<br>always<br>her |
| I will always remember her | her | always<br>remember |

# Improve the Model's Accuracy of Word2Vec

In two ways:

1. Extract the parent nodes of an entity.
2. Better positioning the entity in the walks to maximize its occurrence in the training samples.

For example:

```
("URL#Alice", "URL#knows", "URL#Bob"),
("URL#Alice", "URL#loves", "URL#Bob"),
                    . . .
```

May become:

```
("URL#Casper", "URL#raised", "URL#Alice", "URL#knows", "URL#Bob"),
("URL#Mathilde", "URL#raised", "URL#Alice", "URL#loves", "URL#Bob"),
                    . . .
```

Figure: Word2Vec (Before Improvement).

Figure: Word2Vec (After Improvement).

# FastText

# FastText

Table: Sub-word Generation for Character *N*-Grams of Length 3, 4, 5, and 6.

| Word | Length | Character *n*-grams |
|---|---|---|
| flying | 3 | <fl, fly, lyi, yin, ing, ng> |
| | 4 | <fly, flyi, lyin, ying, ing> |
| | 5 | <flyi, lying, ying> |
| | 6 | <flyin, flying, lying> |

Mainly solves the embeddings creation for Out-Of-Vocabulary words.

# FastText Implementation[11] With KGs

Use the `gensim` library by:

1. Removing the `min_n` and `max_n` parameters for *n*-grams splitting.
2. Allowing to compute *n*-grams for walks only[10] by separating the objects and predicates URIs with a splitting function.
3. Avoiding dependency on Cython.

---

[10]The object nodes in `pyRDF2Vec` are encoded in MD5 to reduce their storage in RAM.

[11] https://github.com/IBCNServices/pyRDF2Vec/blob/feature/bert/pyrdf2vec/embedders/fasttext.py

# Benchmarks

# Benchmarks

- Varied the maximum:
  1. number of walks per entity;
  2. depth per walk.
- Performed on `MUTAG`.
- Use 320 training entities and attempt to predict 68 test entities.
- Take the mean and standard deviation of five tests.

Embedding techniques and walking strategies:

- IDLab's servers with 4 CPUs, 64 GB RAM, and one GPU.

Sampling strategies:

- ThinkPad machine with 4 CPUs, 16 GB RAM.

# Benchmarks for Embedding Techniques

Figure: Embedding Techniques for `MUTAG` (Maximum Depth per Walk).

# Benchmarks for Embedding Techniques

Figure: Embedding Techniques for `MUTAG` (Maximum Number of Walks).

# Strategies

# SplitWalker[12]

- Extension of the `RandomWalker`.
- Based on the principle of preprocessing.
- Customizable by a splitting function.
- Inspired by the way FastText works.

[12] https://github.com/IBCNServices/pyRDF2Vec/blob/master/pyrdf2vec/walkers/split.py

# SplitWalker

| | Initial Node | Node After Splitting |
|---|---|---|
| Walk 1 | http://dl-learner.org/carcinogenesis#d19 | http://dl-learner.org/carcinogenesis#d19 |
| | http://dl-learner.org/carcinogenesis#hasBond | has |
| | | bond |
| | http://dl-learner.org/carcinogenesis#bond3209 | 3209 |
| Walk 2 | http://dl-learner.org/carcinogenesis#d36 | http://dl-learner.org/carcinogenesis#d36 |
| | http://www.w3.org/1999/02/22-rdf-syntax-ns#type | type |
| | http://dl-learner.org/carcinogenesis#bond | bond |

Table: Walk Transformation With `SplitWalker`.

# Benchmarks for Walking Strategies

| Walker | Average Rank |
|---|:---:|
| HALKWalker(freq_threshold=0.01) | 1 |
| SplitWalker | 2 |
| RandomWalker | 3 |
| NGramWalker(grams=3) | 4 |
| WalkletWalker | 5 |
| AnonymousWalker | 6 |

Table: Average Rank of the Walking Strategies.

# WideSampler[13]

- Maximizes the extraction of shared features between entities.
- Assumes that single entity-specific features have a negligible impact on the quality of the generated embeddings.
- Assigns higher weights to edges that lead to the largest number of:
  1. predicates and objects in the neighborhood;
  2. occurrence of predicates and objects in a graph.

---

[13] https://github.com/IBCNServices/pyRDF2Vec/blob/master/pyrdf2vec/samplers/wide.py

# Benchmarks for Sampling Strategies

| Sampler | Average Rank |
|---|---|
| `WideSampler` | 1 |
| `ObjPredFreqSampler(inverse=True)` | 2 |
| `PredFreqSampler` | 3 |
| `ObjFreqSampler` | 4 |
| `ObjFreqSampler(inverse=True)` | 5 |
| `PageRankSampler(inverse=True,alpha=0.85)` | 6 |
| `PageRankSampler(inverse=True,split=True,alpha=0.85)` | 7 |
| `ObjFreqSampler(inverse=True,split=True)` | 8 |
| `PageRankSampler(split=True,alpha=0.85)` | 9 |
| `PredFreqSampler(inverse=True)` | 10 |
| `ObjPredFreqSampler` | 11 |
| `UniformSampler` | 11 |
| `PageRankSampler(alpha=0.85)` | 12 |

Table: Average Rank of the Sampling Strategies.

# Future Work

# Future Work

- Inject `(subject,object)` pairs to BERT instead of pairs of walks.
- Evaluate:
  - the KG-oriented BERT models with other embedding techniques;
  - BERT, `SplitWalker` and `WideSampler` on larger KGs.
- If necessary, implement a new KG-oriented BERT model.

# Conclusion

# Conclusion

- The traditional BERT model has not been as effective as expected: small gain for the model's accuracy in counterpart of a more important training time and storage space.
- `WideSampler` and `SplitWalker` finished first and second respectively in average rank of the `MUTAG` benchmarks.
- Improved the accuracy of the Word2Vec model.
- Gain in popularity for `pyRDF2Vec`.

Question Time 🐈

# Appendices

# Related Work

Three categories of conversion algorithms based:

1. Direct encoding (e.g., TransE[14]);



Figure: Embedding Computations With TransE.

2. Deep Learning (e.g., R-GCNs[15]);
3. Path/Walk (e.g., RDF2Vec[16]).

---

[14]Translating Embeddings.
[15]Modeling Relational Data with Graph Convolutional Networks.
[16]Resource Description Framework To Vector.

# Continuous Bag-of-Words Model



Figure: CBOW Model Architecture.

# Skip-Gram Model



Figure: Skip-Gram Model Architecture.

# Recurrent Neural Networks



Figure: Sequence-to-Sequence Learning With RNNs.

# Attention Mechanism



Figure: Seq2Seq Learning With Attention Mechanism.

# Transfomer

Figure: Transformer Model Architecture.

Source: VASWANI et al. – Attention Is All You Need

# Input Embeddings



Figure: Example of Sentence Pair Encoding.

Source: DEVLIN et al. – BERT

# Masked Language Model

Table: Example of MLM With BERT.

best
↑
GNU   Emacs   is   the   [MASK]   text   editor

Each input sequence usually has 15 % of its tokens hidden where:

- a token is replaced by a [MASK] token in 80 % of the cases;
- a token is replaced randomly in 10 % of the cases;
- a token remains unchanged in 10 % of the cases.

# Next Sentence Prediction

Table: Example of NSP With BERT.

| Sentence A | Sentence B | Label |
|---|---|---|
| GNU Emacs is the best text editor. | I should use it. | `isNextSentence` |
| GNU Emacs is the best text editor. | I have a cat. | `NotNextSentence` |

---

**Algorithm** `get_weight(h, d, c)`

---

**Require:** a $\mathcal{H}$ 2-tuple that contains a predicate and an object.
**Require:** a $\mathcal{D}$ array of $n \geq 1$ degree indexed from 0 to $n - 1$
**Require:** a $\mathcal{C}$ array of $n \geq 1$ counter of neighbors indexed from 0 to $n - 1$
**Ensure:** The weight of the hop for this predicate
  1: **if** $\mathcal{D}_{preds}$ and $\mathcal{D}_{objs}$ and $\mathcal{C}$ **then**
  2:    **return** $(\mathcal{C}[\mathcal{H}[0]_{name}] + \mathcal{C}[\mathcal{H}[1]_{name}])$
$$\left( \frac{\mathcal{D}_{preds}[\mathcal{H}[0]_{name}] + \mathcal{D}_{objs}[\mathcal{H}[1]_{name}]}{2} \right)$$
  3: **end if**

**Listing:** Splits Nodes of Random Walks with `SplitWalker` (Part I).

```python
def basic_split(self, walks):
    canonical_walks = set()
    for walk in walks:
        canonical_walk = [walk[0].name]
        for i, _ in enumerate(walk[1::], 1):
            vertices = []
            if "http" in walk[i].name:
                vertices = " ".join(re.split("[\#]", walk[i].name)).split()
            if i % 2 == 1:
                name = vertices[1] if vertices else walk[i].name
                preds = [
                    sub_name
                    for sub_name in re.split(r"([A-Z][a-z]*)", name)
                    if sub_name
                ]
                for pred in preds:
                    canonical_walk += [pred.lower()]
```

```
...
    else:
      name = vertices[-1] if vertices else walk[i].name
      objs = []
      try:
        objs = [str(float(name))]
      except ValueError:
        objs = re.sub("[^A-Za-z0-9]+", " ", name).split()
        if len(objs) == 1:
          match = re.match(
            r"([a-z]+)([0-9]+)", objs[0], re.I
          )
          if match:
            objs = list(match.groups())
      for obj in objs:
        canonical_walk += [obj.lower()]
  canonical_walk = list(dict(zip(canonical_walk, canonical_walk)))
  canonical_walks.add(tuple(canonical_walk))
  return canonical_walks
```