



AUTHORSHIP VERIFICATION

GRAPH-BASED SIAMESE NETWORK FOR AUTHORSHIP VERIFICATION

Daniel Embarcadero Ruiz

Dra Helena Gómez Adorno

Ivan Reyes Hernández

Alexis Garcia

Alberto Embarcadero

UNAM, MEXICO

OUTLINE

- Authorship Verification
- Model Overview
- Relevant details:
 - ◆ Modeling texts as graphs
 - ◆ Graph-based Siamese Network
 - ◆ GBSN Ensemble
- Submitted model
- Results
- Conclusions

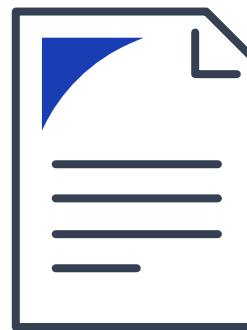
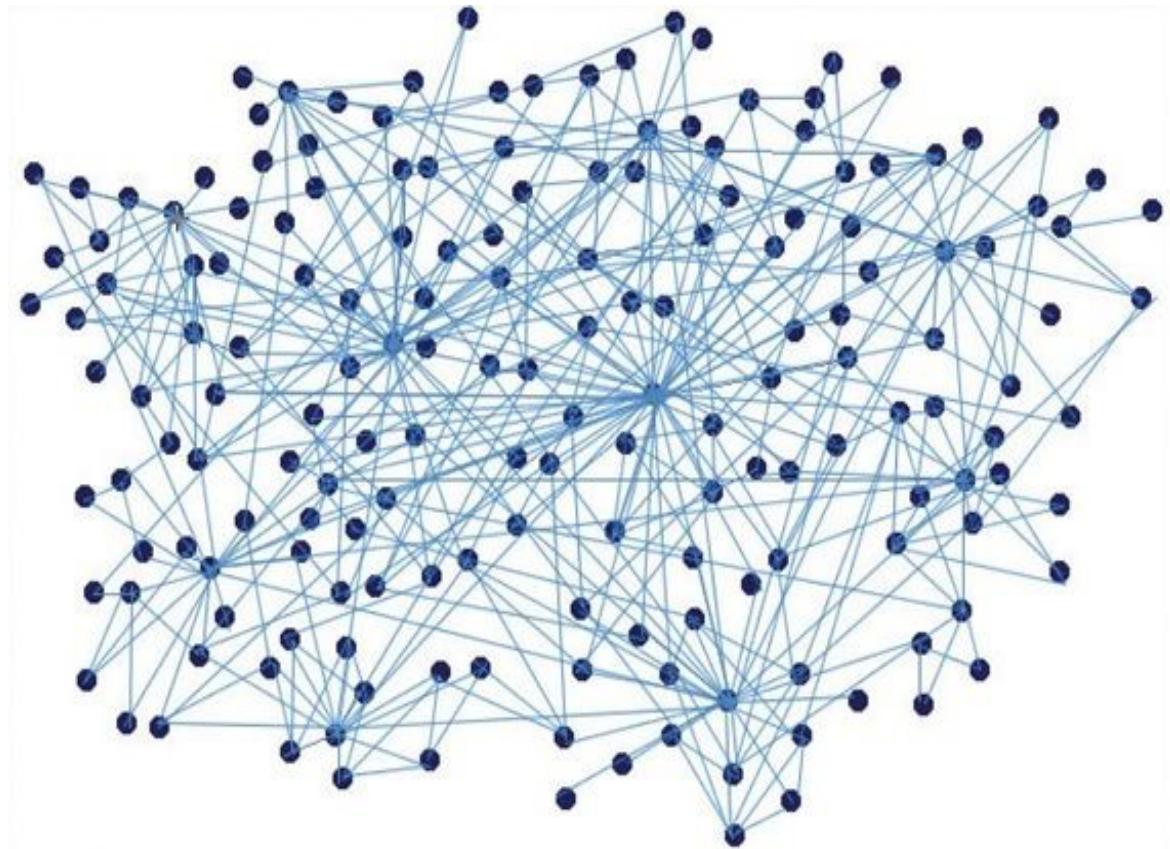
AUTHORSHIP VERIFICATION

- Is the task of determining if two texts were written by the same author.



MODEL OVERVIEW

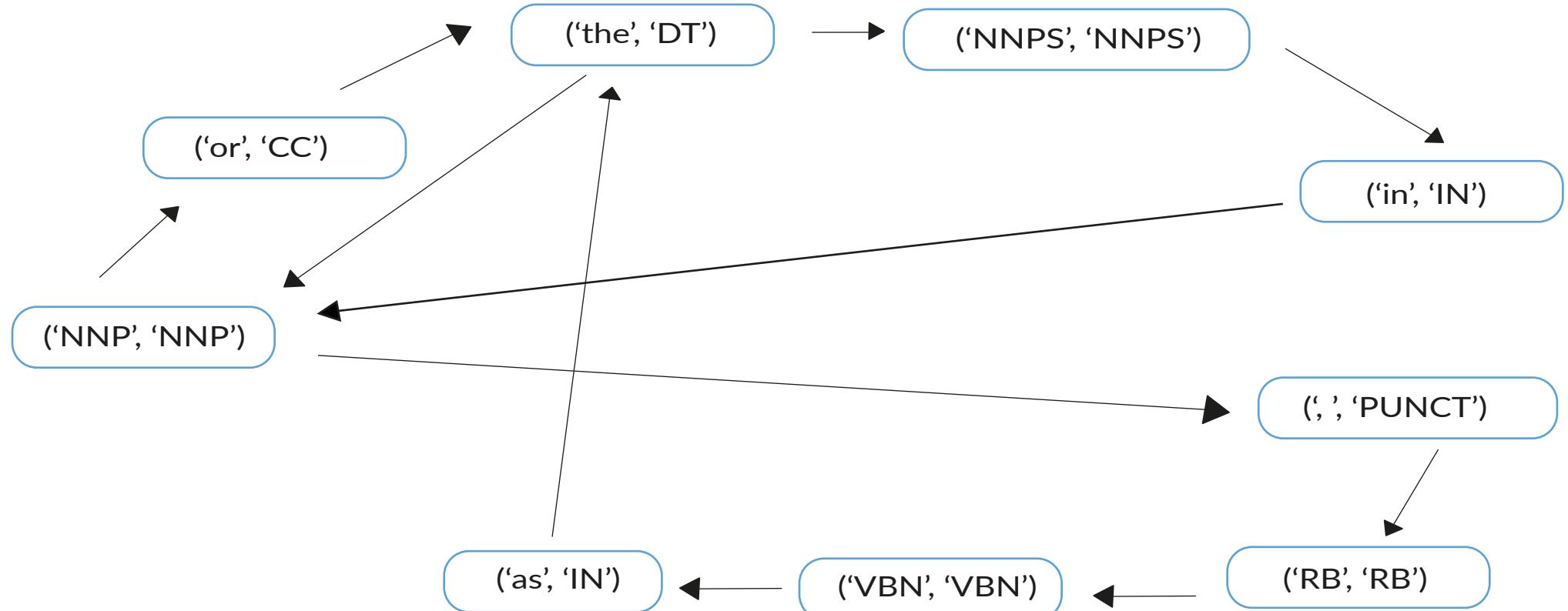
- Represent texts as graphs.
- Let a graph neural network automatically extract features.
- Use these ideas inside a Siamese Network.



MODELING TEXTS AS GRAPHS

Input: “Momo, also known as The Grey Gentlemen or the Men in Grey”

REDUCE_LABELS = {NNP, NNPS, VBN, RB}



PREPROCESS

Input: “Momo, also known as The Grey Gentlemen or the Men in Grey”

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

- Substitute non-ASCII characters.
- Tokenize.
- Get POS labels (we use Treebank plus PUNCT and OTHER).
- Normalize to lowercase.

PREPROCESS

Input: "Momo, also known as The Grey Gentlemen or the Men in Grey"

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed:

```
[ ('momo', 'NNP'), (',', 'PUNCT'), ('also', 'RB'),  
  ('known', 'VBN'), ('as', 'IN'), ('the', 'DT'),  
  ('grey', 'NNP'), ('gentlemen', 'NNP'),  
  ('or', 'CC'), ('ther', 'DT'), ('men', 'NNPS'),  
  ('in', 'IN'), ('grey', 'NNP') ]
```

GRAPH CONSTRUCTION

Input: “Momo, also known as The Grey Gentlemen or the Men in Grey”

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', ',', 'PUNCT'), ...]

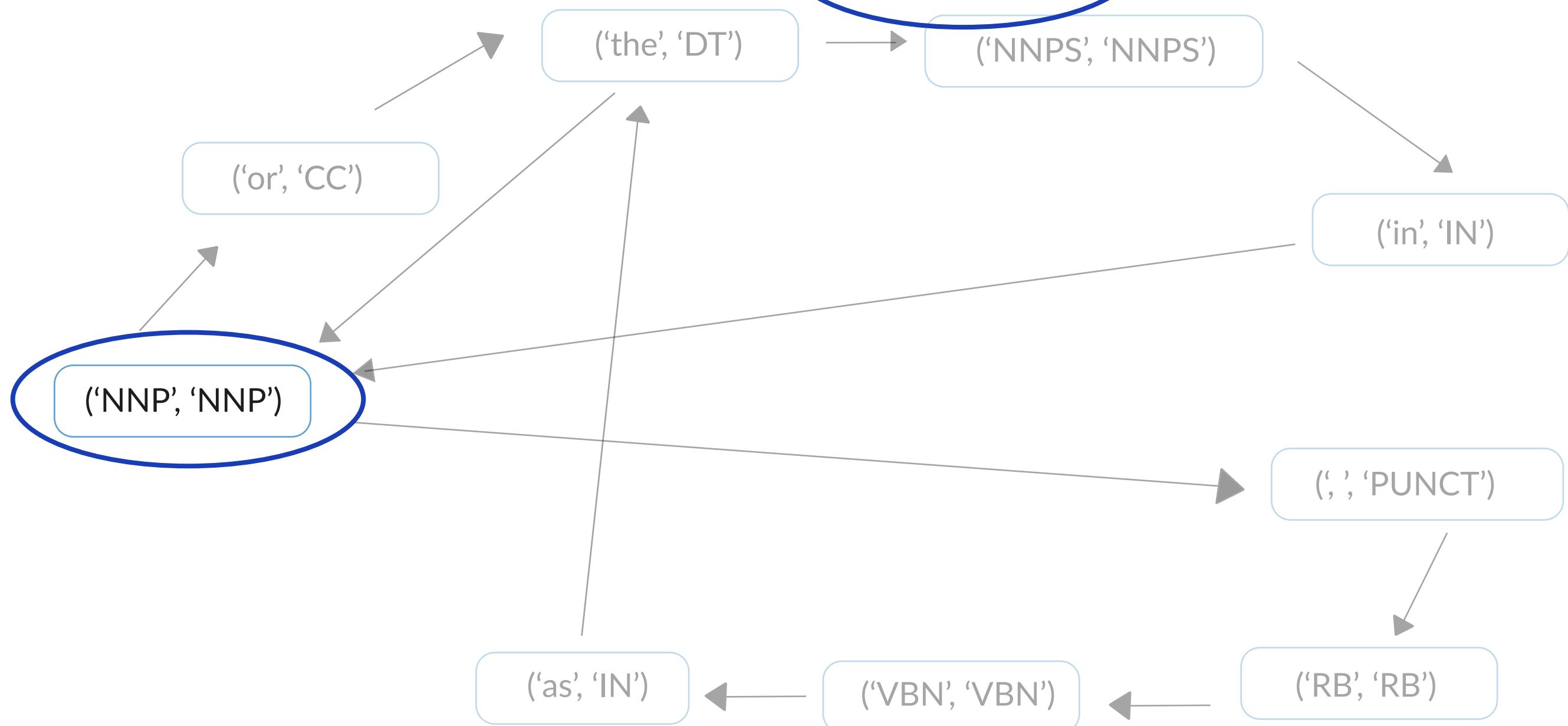
- It is a co-occurrence based construction.
- The nodes will be tuples (token, label).
- The edges occur if and only if the tokens appear together in text.

GRAPH CONSTRUCTION

Input: "Momo, also known as The Grey Gentlemen or the Men in Grey"

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', ',', 'PUNCT'), ...]

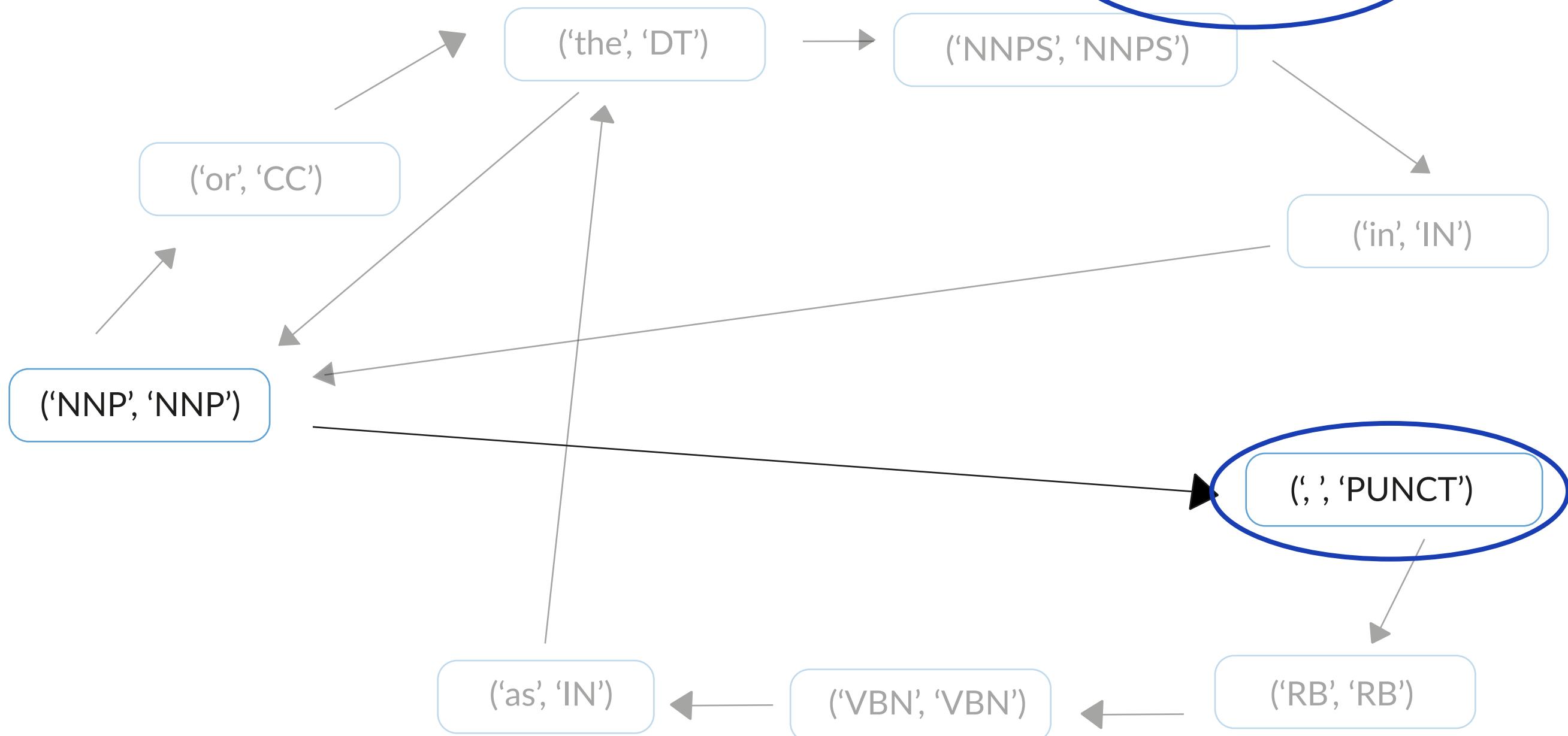


GRAPH CONSTRUCTION

Input: "Momo, also known as The Grey Gentlemen or the Men in Grey"

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', ',', 'PUNCT'), ...]

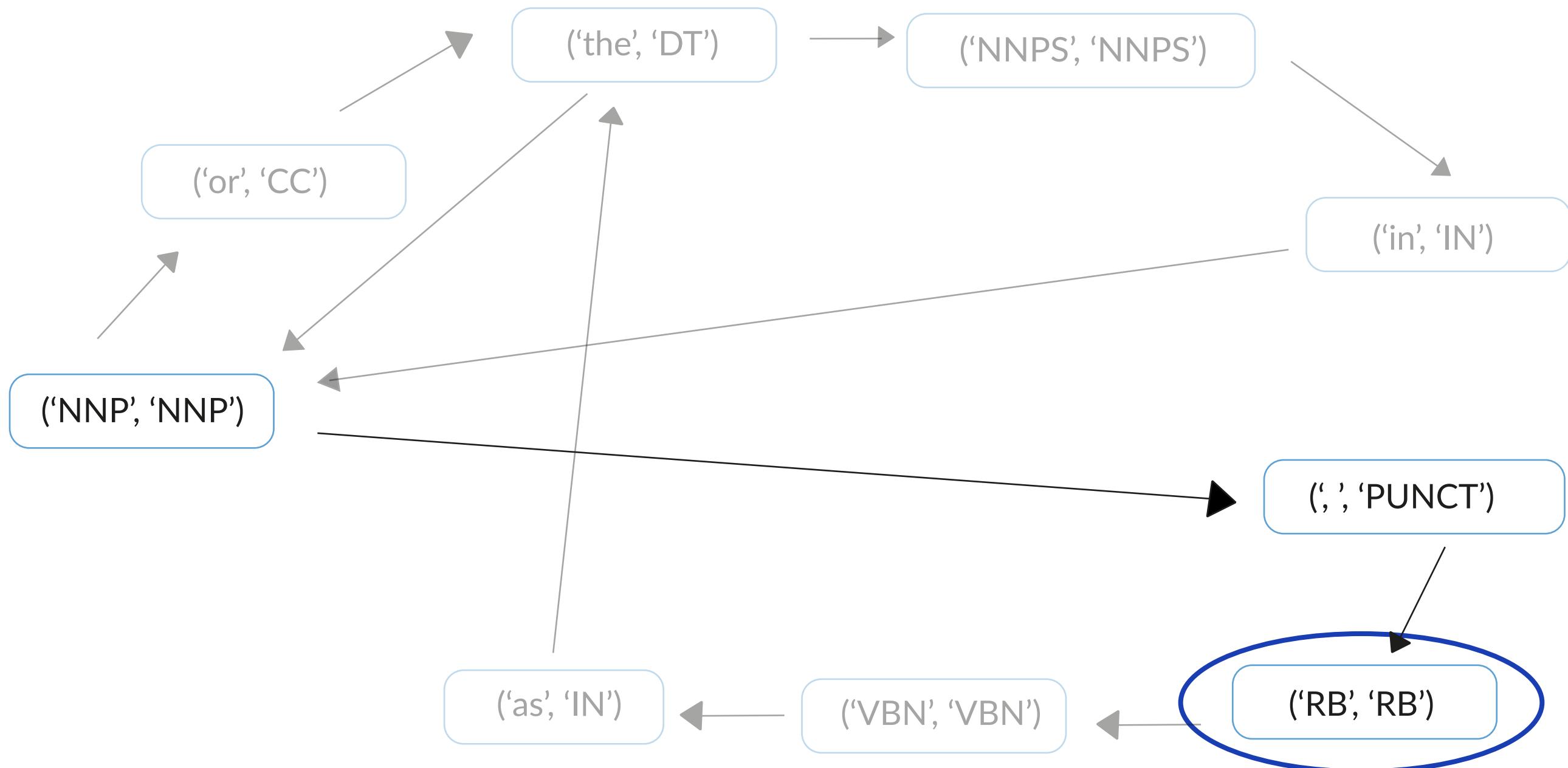


GRAPH CONSTRUCTION

Input: "Momo, also known as The Grey Gentlemen or the Men in Grey"v

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', 'PUNCT'), ...]

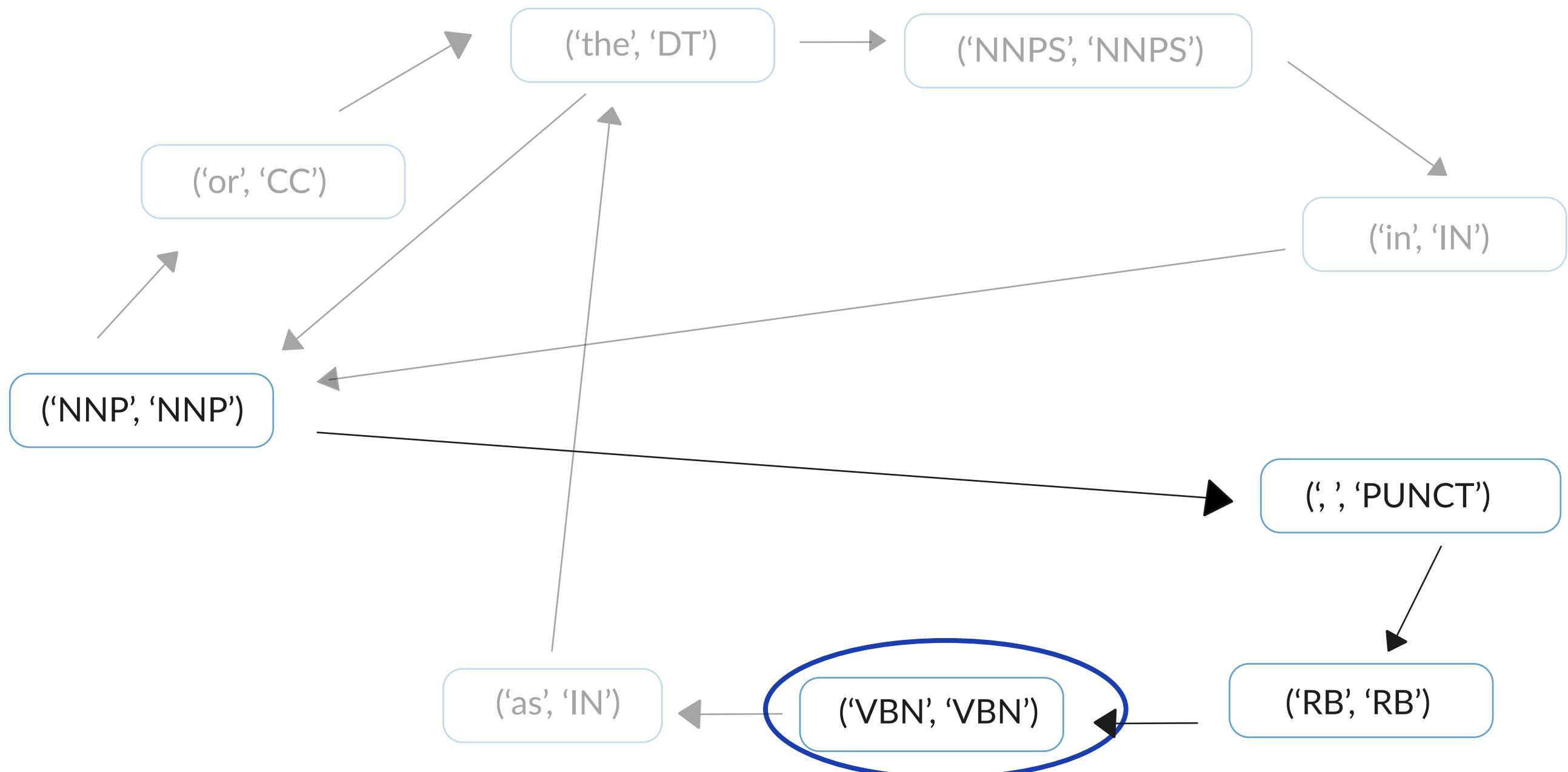


GRAPH CONSTRUCTION

Input: "Momo, also known as The Grey Gentlemen or the Men in Grey"

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', ',', 'PUNCT'), ...]

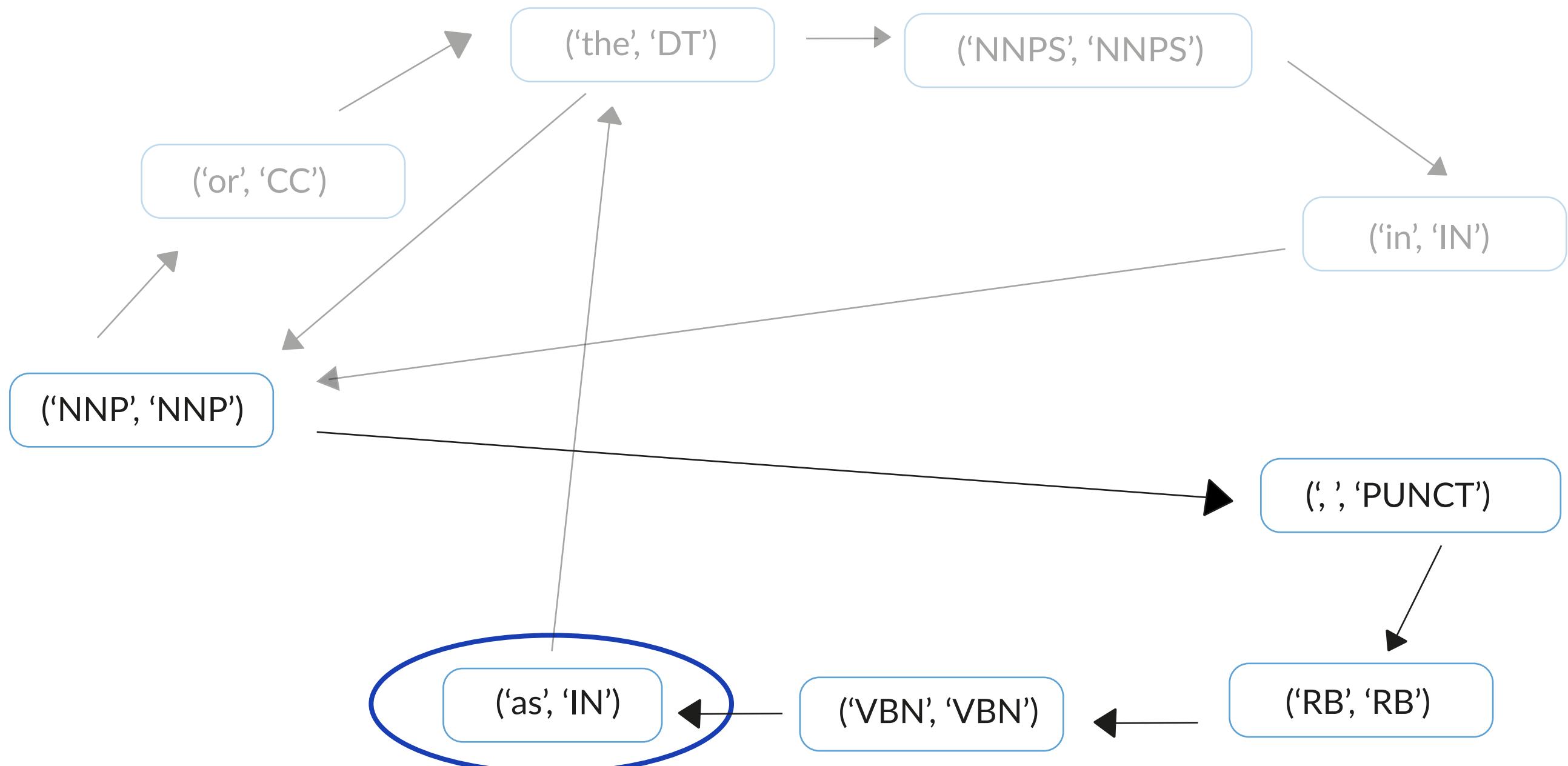


GRAPH CONSTRUCTION

Input: "Momo, also known as The Grey Gentlemen or the Men in Grey"

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', ',', 'PUNCT'), ...]

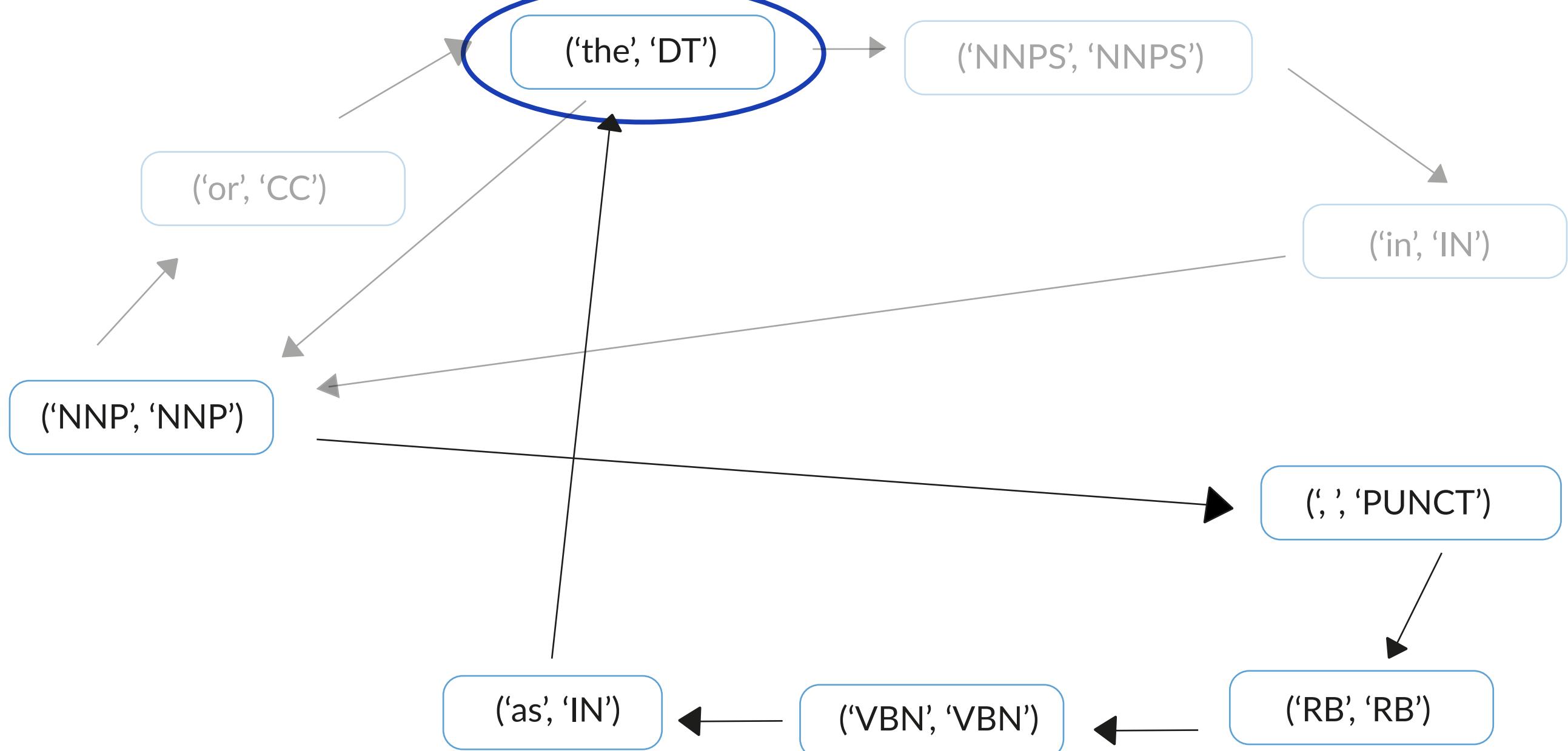


GRAPH CONSTRUCTION

Input: "Momo, also known as The Grey Gentlemen or the Men in Grey"

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', ',', 'PUNCT'), ...]

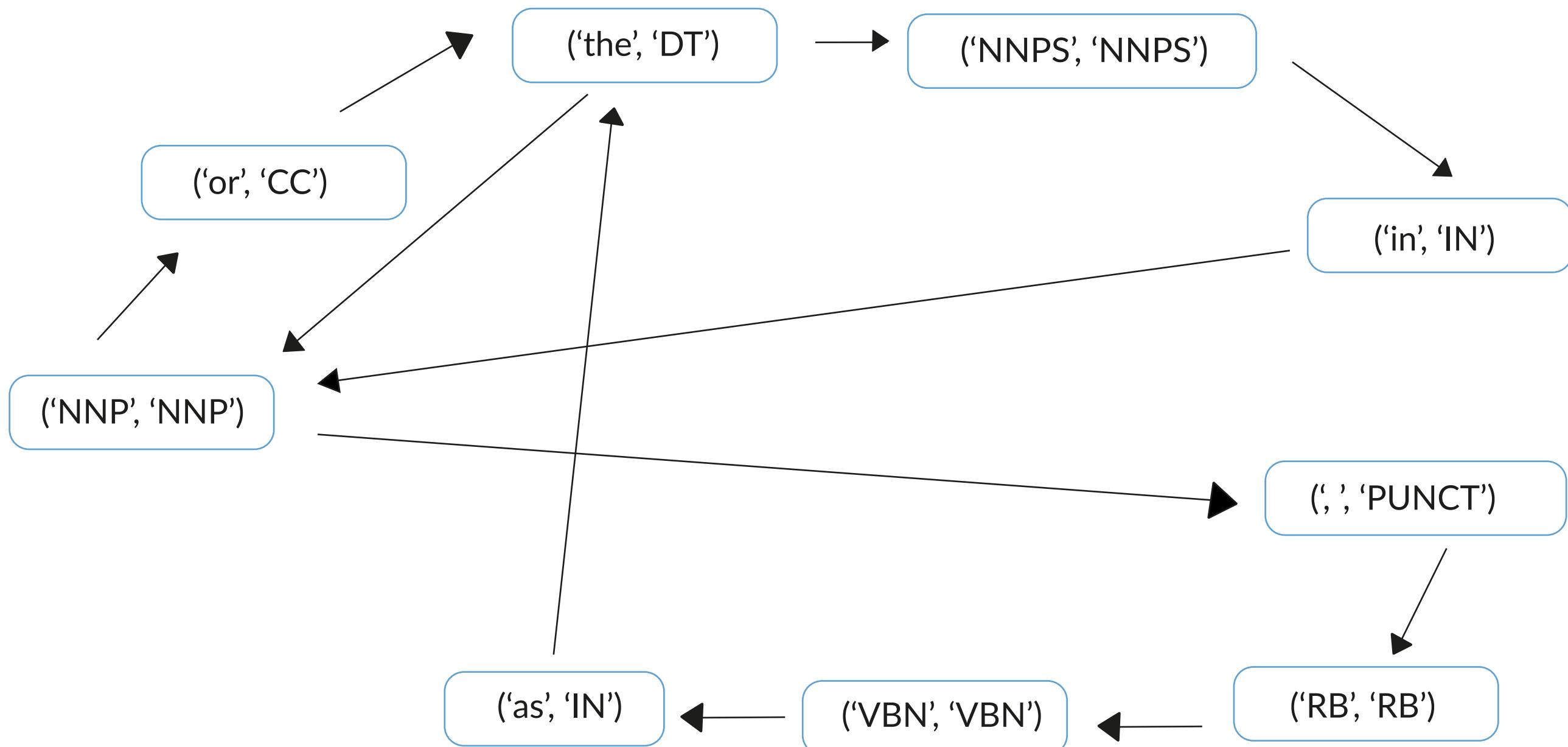


GRAPH CONSTRUCTION

Input: "Momo, also known as The Grey Gentlemen or the Men in Grey"

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', 'PUNCT'), ...]

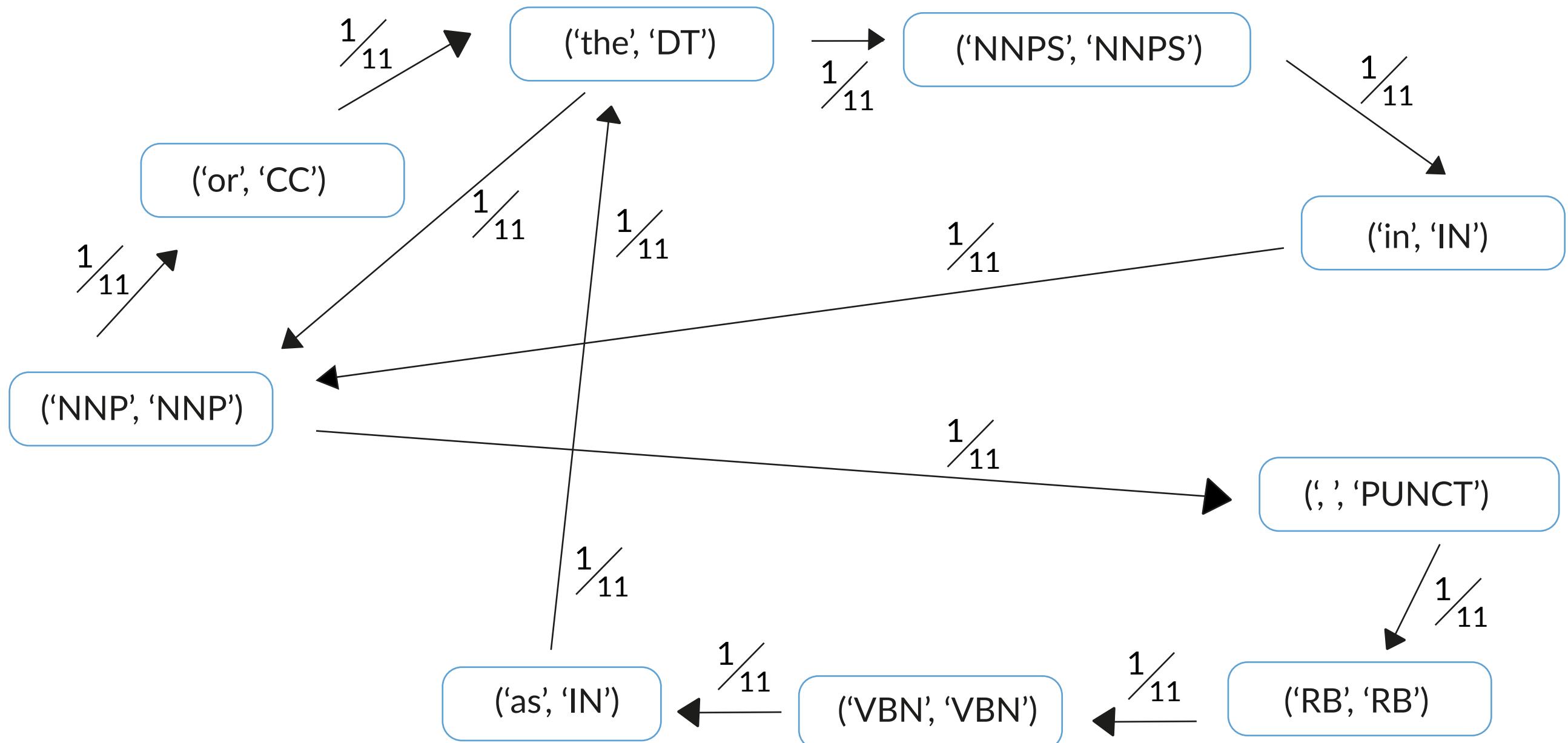


GRAPH CONSTRUCTION

Input: “Momo, also known as The Grey Gentlemen or the Men in Grey”

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

Processed: [('Momo', 'NNP'), (',', ',', 'PUNCT'), ...]

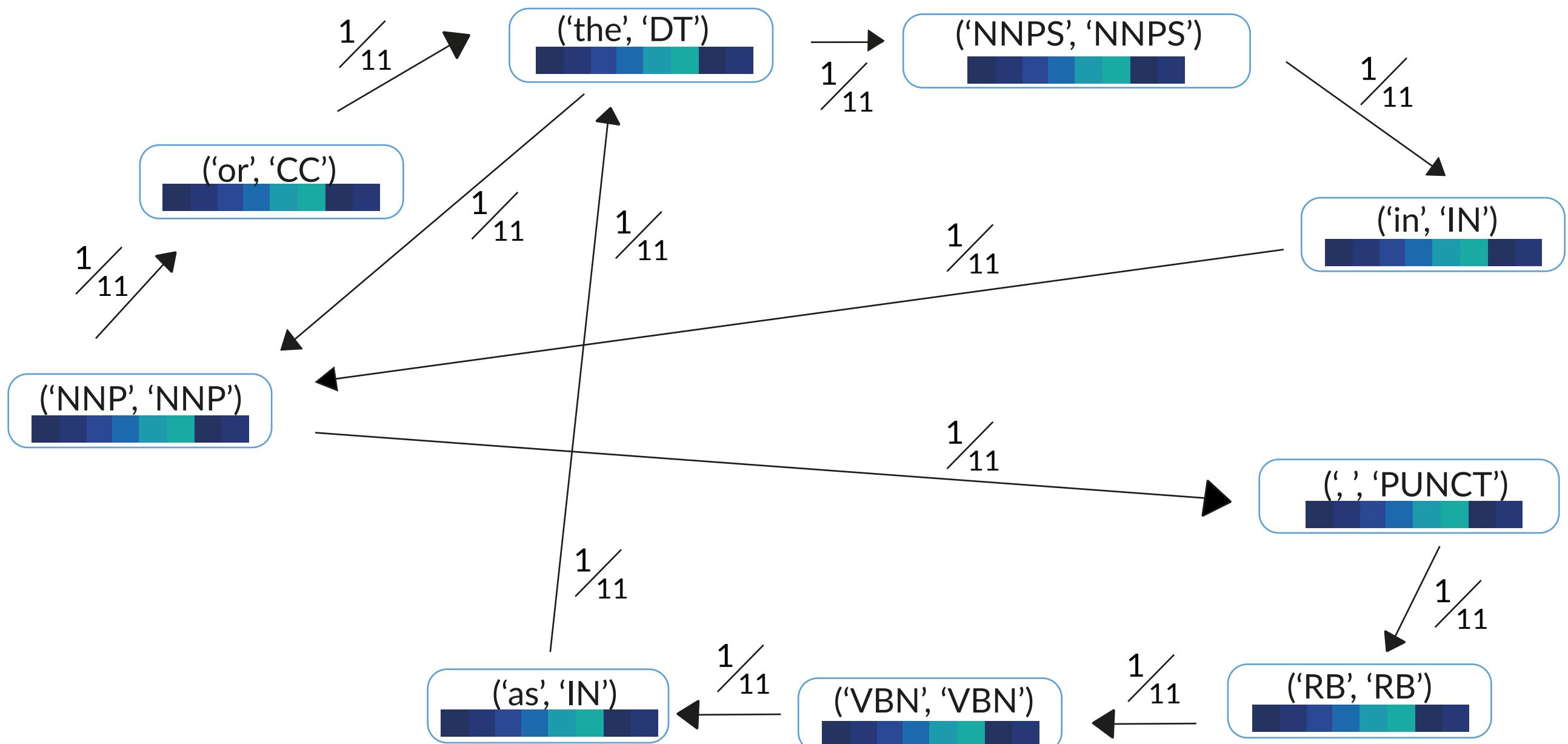


GRAPH CONSTRUCTION

Input: “Momo, also known as The Grey Gentlemen or the Men in Grey”

REDUCE_LABELS = {NNP, NNPS, VBN, RB}

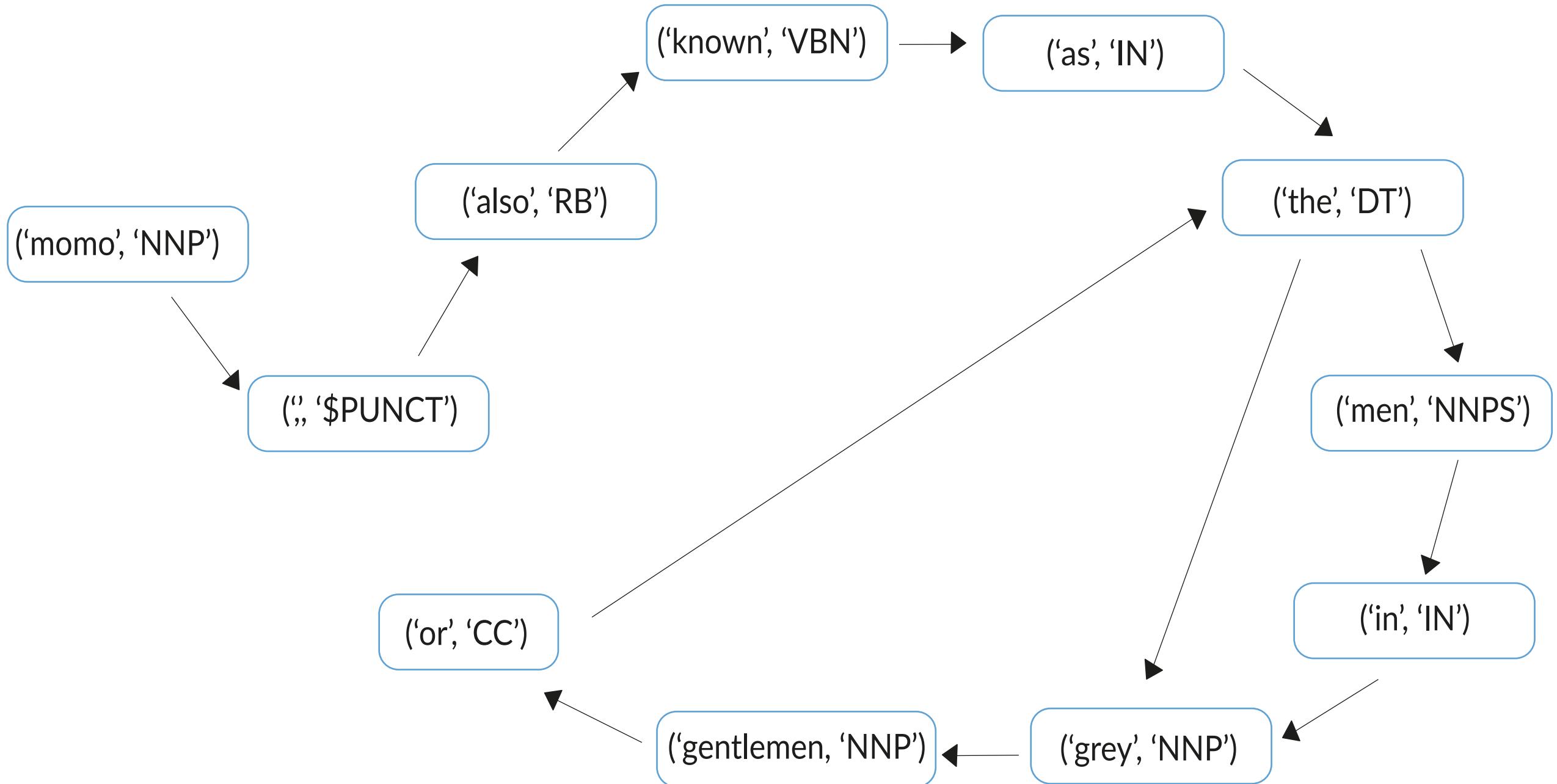
Processed: [('Momo', 'NNP'), (', ', ', ', 'PUNCT'), ...]



Input: "Momo, also known as The Grey Gentlemen
or the Men in Grey

FULL GRAPH (COOCURRENCE)

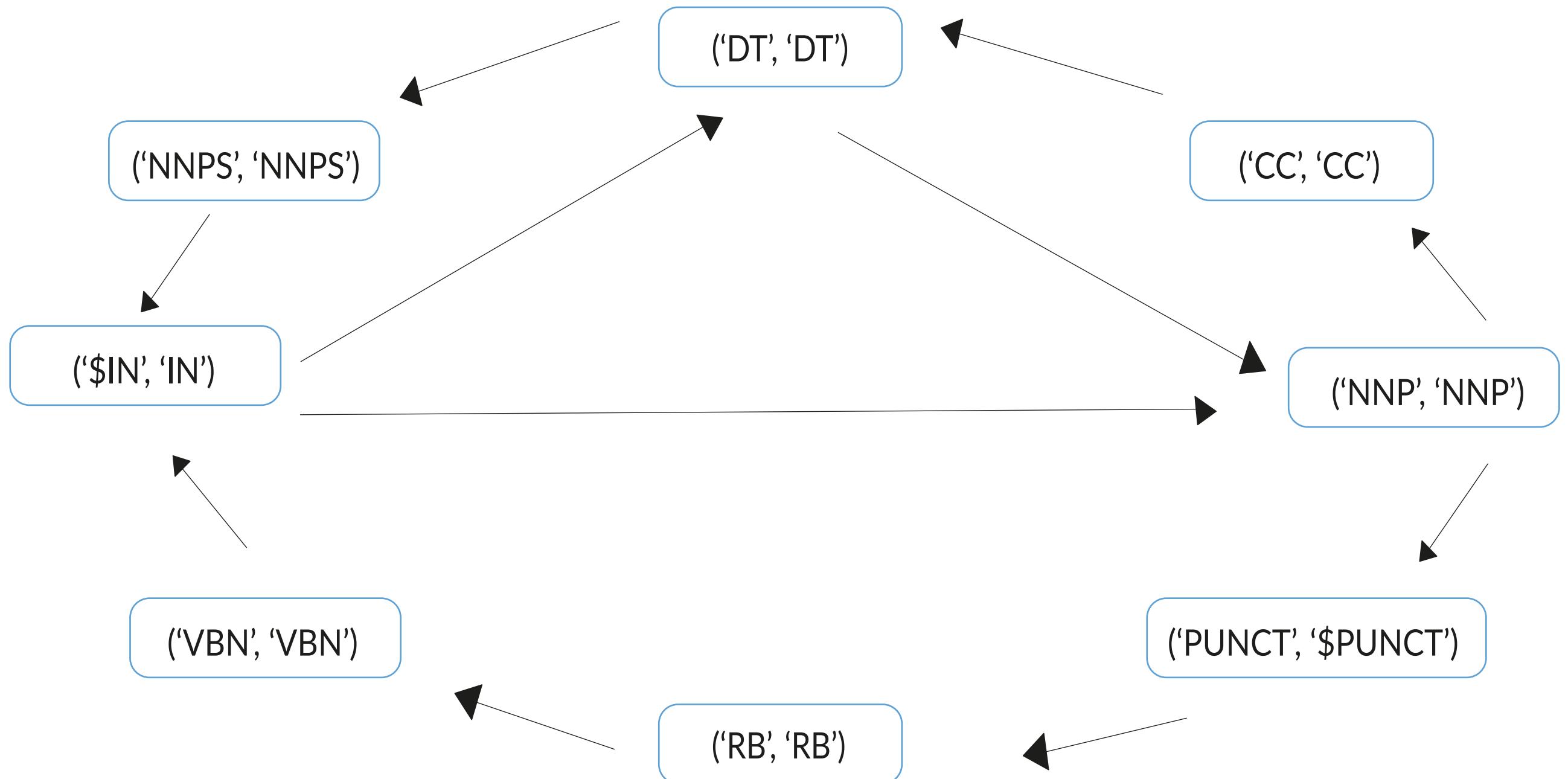
REDUCE_LABELS = { }



Input: "Momo, also known as The Grey Gentlemen
or the Men in Grey

SHORT GRAPH

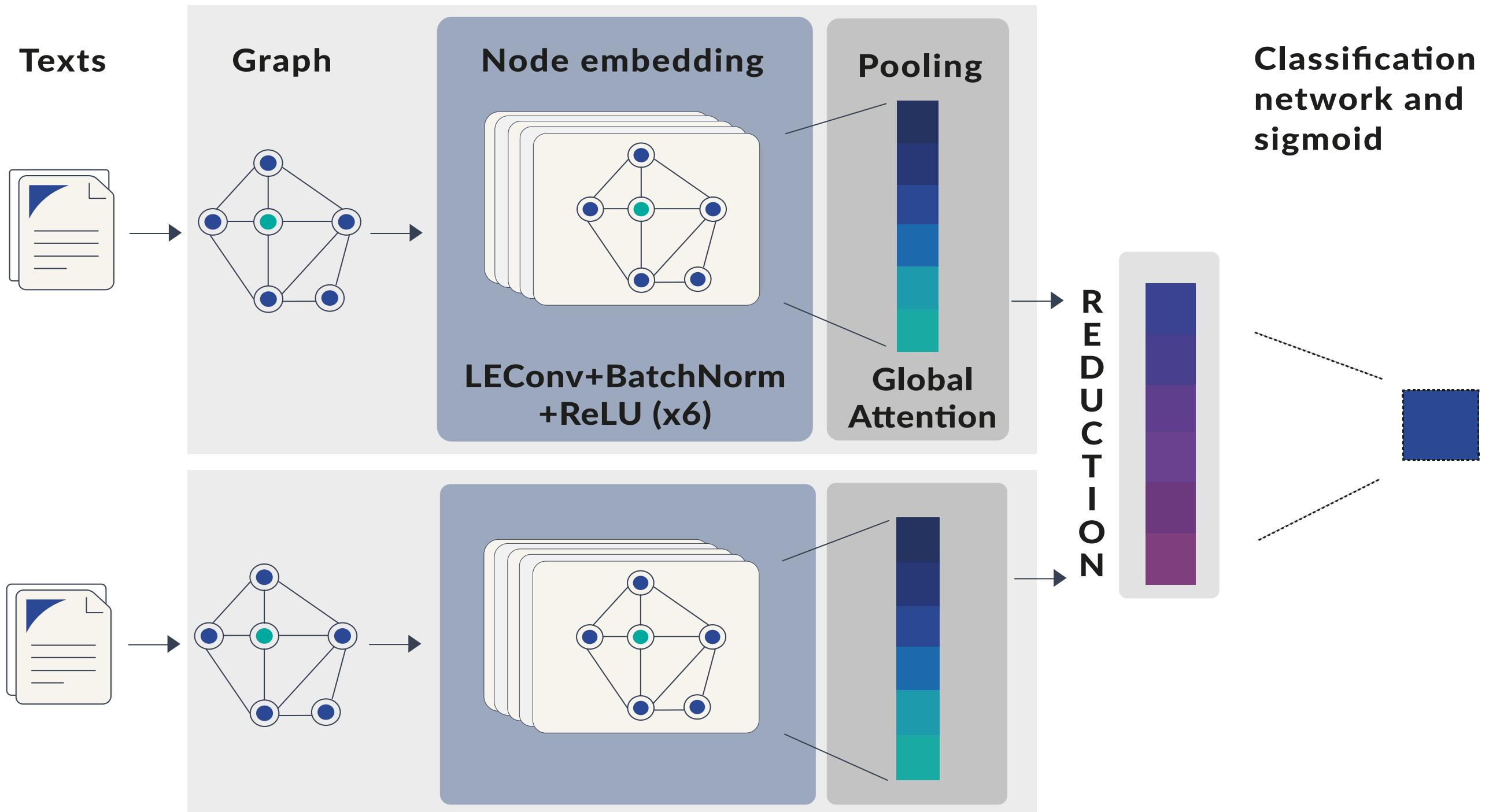
REDUCE_LABELS = ALL LABELS



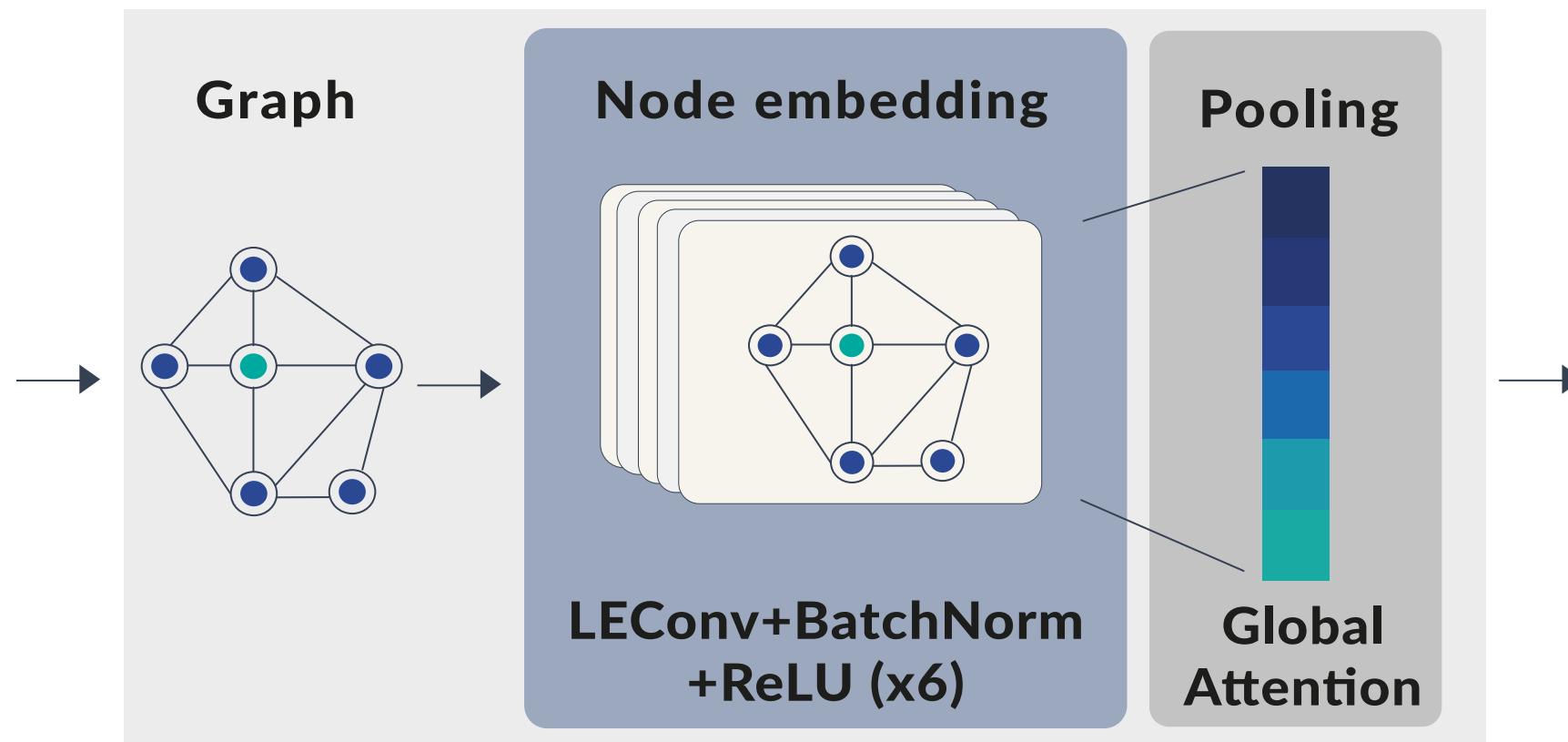
GRAPH SIZE

Nombre	REDUCE_LABELS set	Nodes	Edges
Short	All labels	33	407
Med	Labels corresponding to adjectives, nouns, adverbs, verbs, cardinal numbers, foreign words, list item marker, symbols and others	138	1168
Full	Empty set	1207	3424

GRAPH-BASED SIAMESE NETWORK



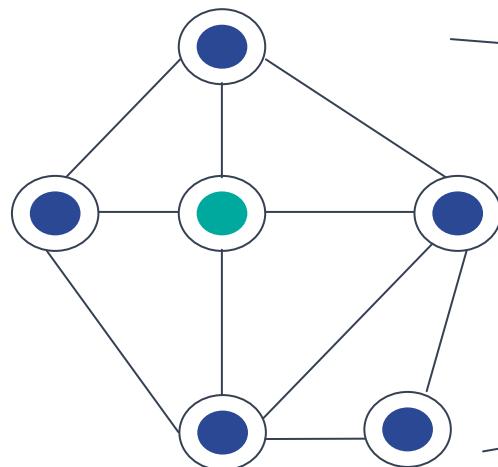
FEATURE EXTRACTION COMPONENTS



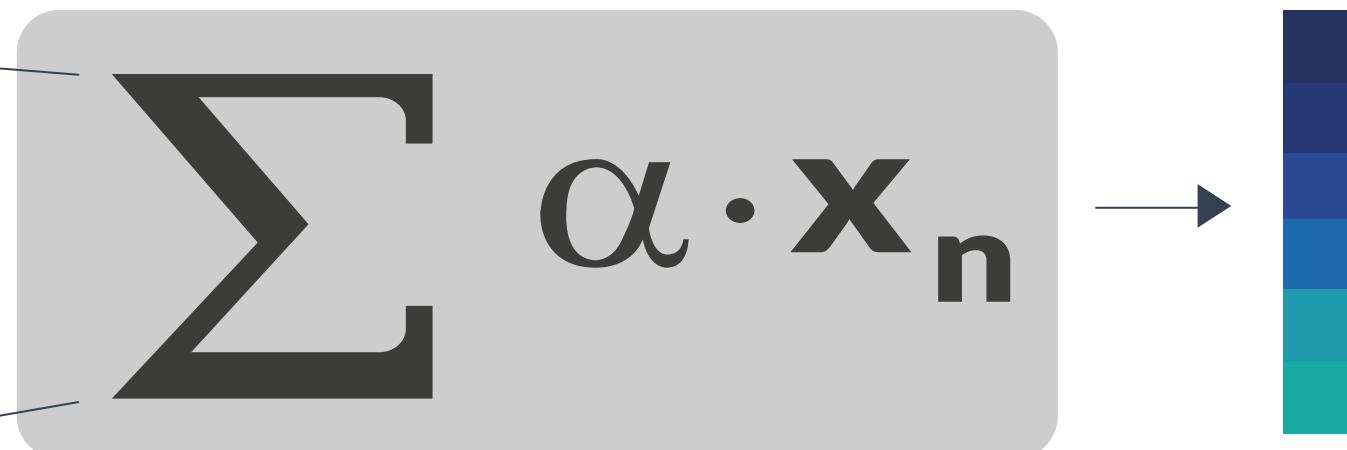
* Local Extreme Convolution (LEConv) layer was originally proposed in:
E. Ranjan, et al. "ASAP: Adaptive Structure Aware Pooling for Learning Hierarchical Graph Representations". Thirty-Fourth AAAI Conference on Artificial Intelligence, Feb 2020.

GLOBAL POOLING

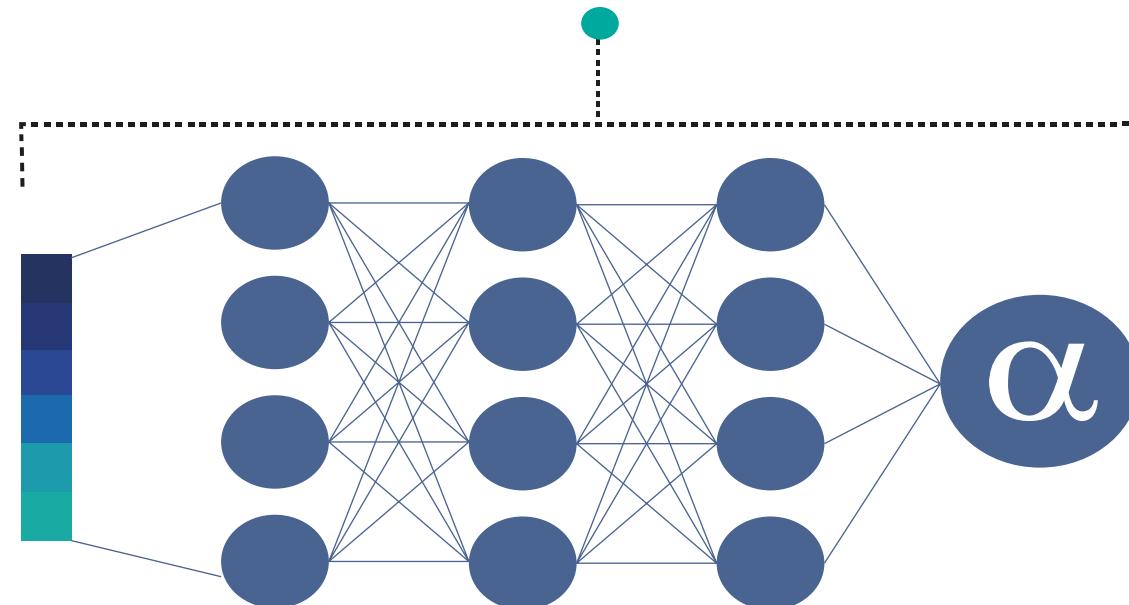
Last layer output



Global attention



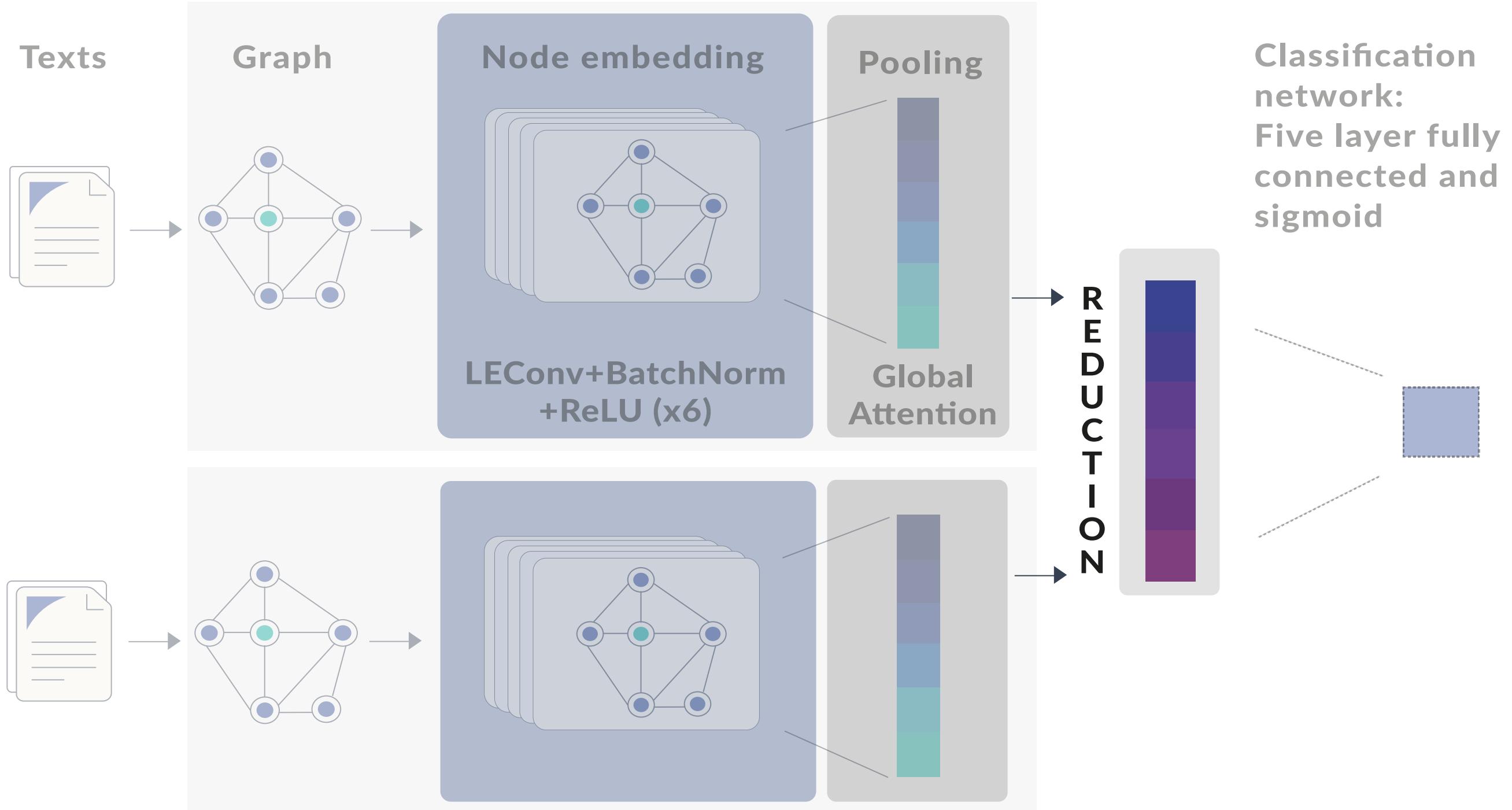
L-pool
Layers



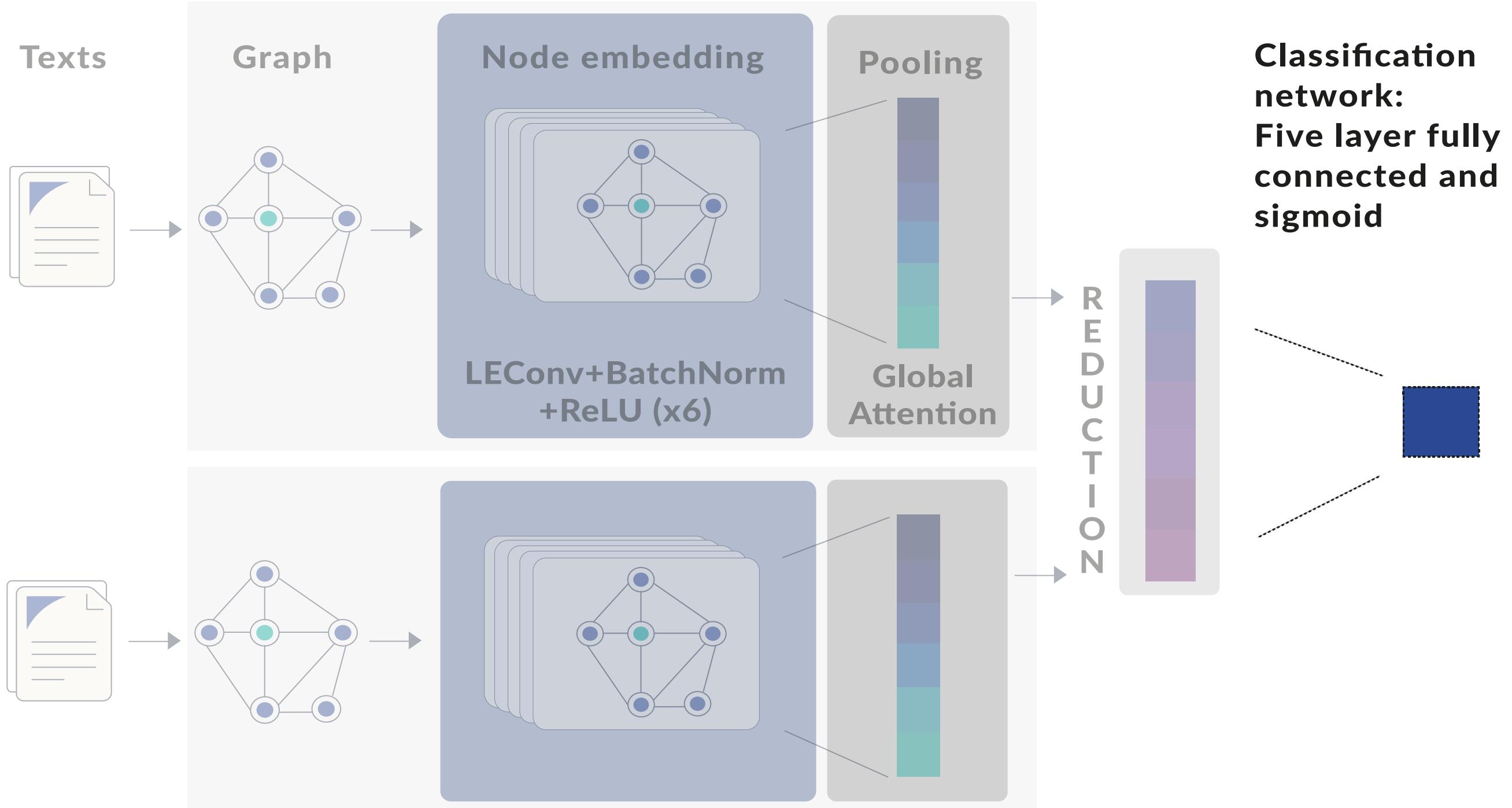
* Global pooling layer was original in:

Y. Li, et. al. "Gated Graph Sequence Neural Networks", Proceedings of ICLR, Sept. 2017.

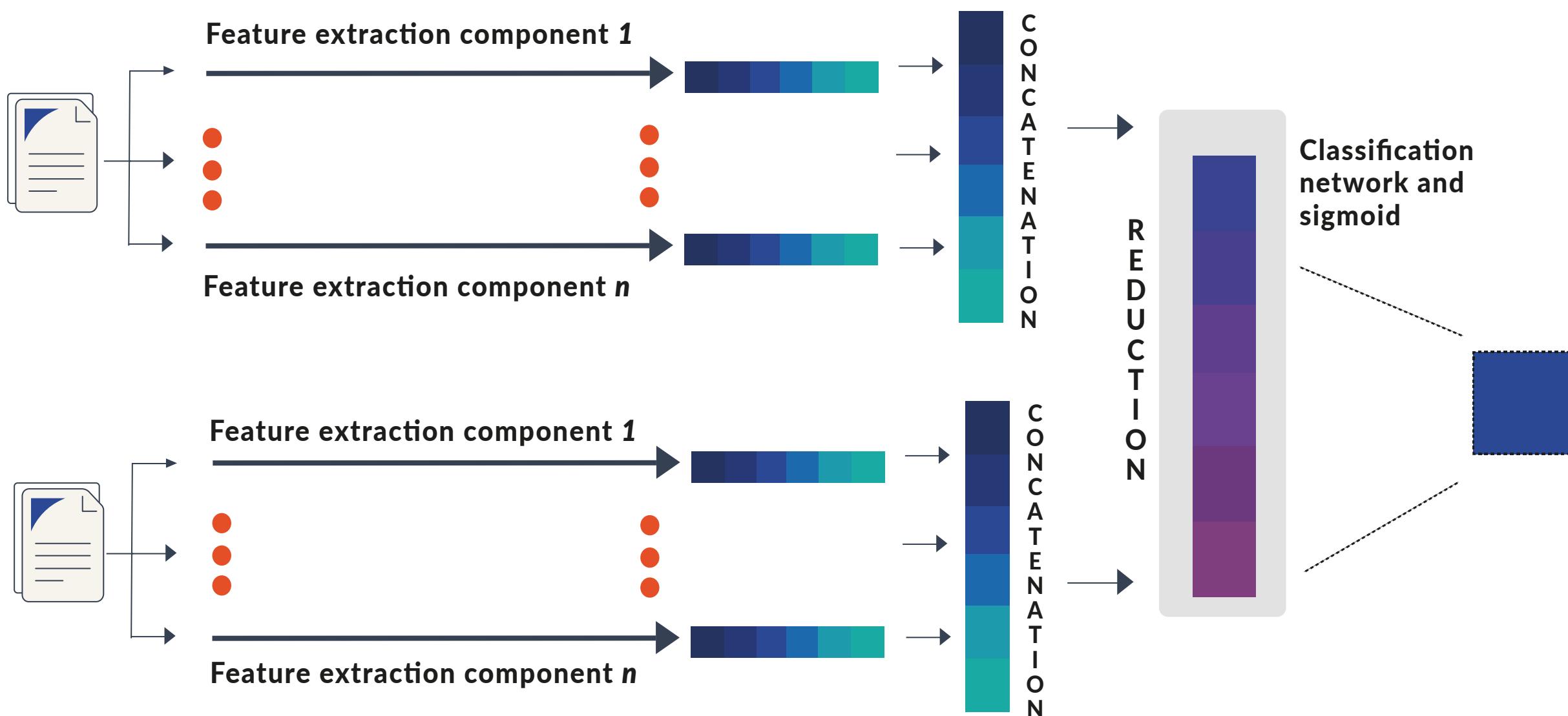
REDUCTION AND CLASSIFICATION



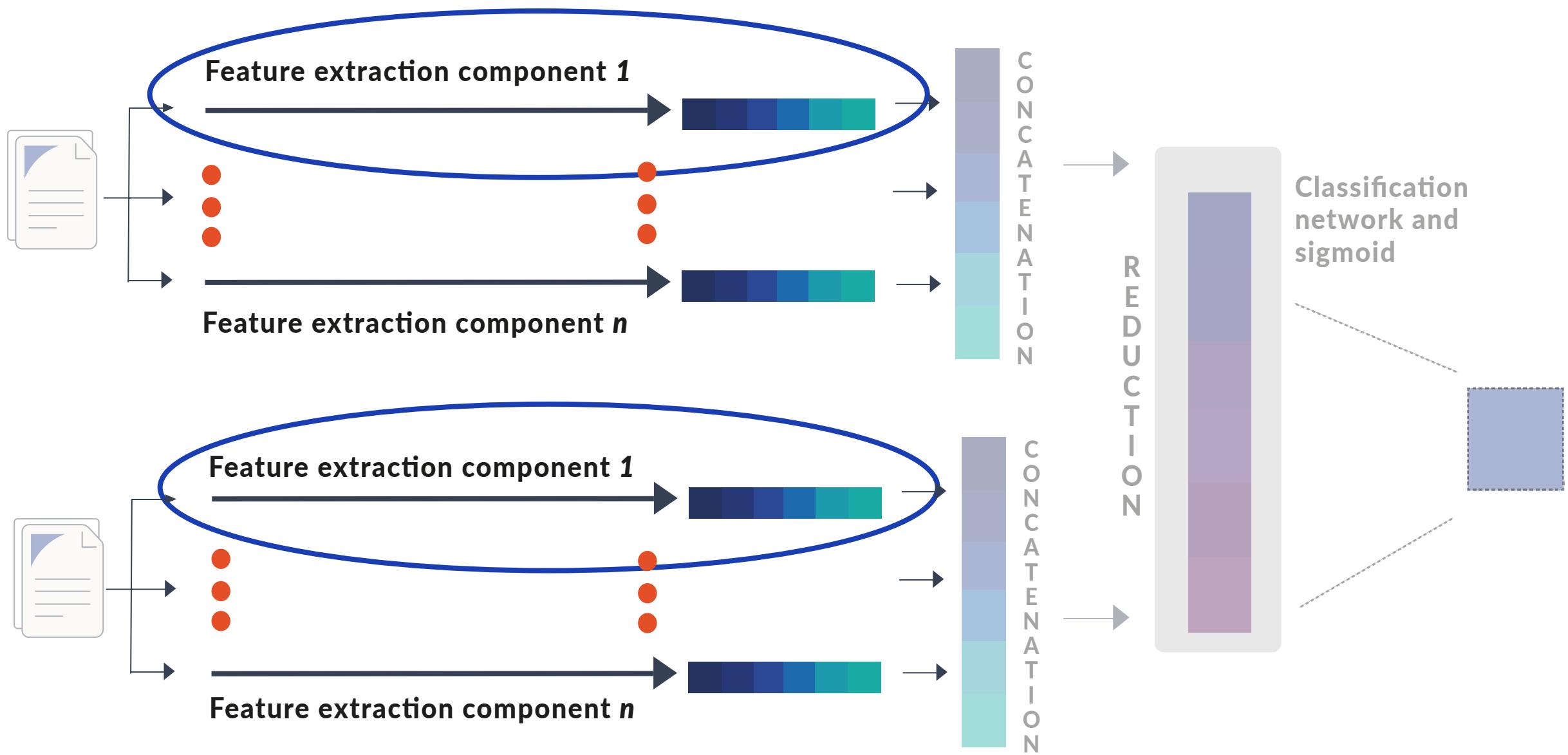
REDUCTION AND CLASSIFICATION



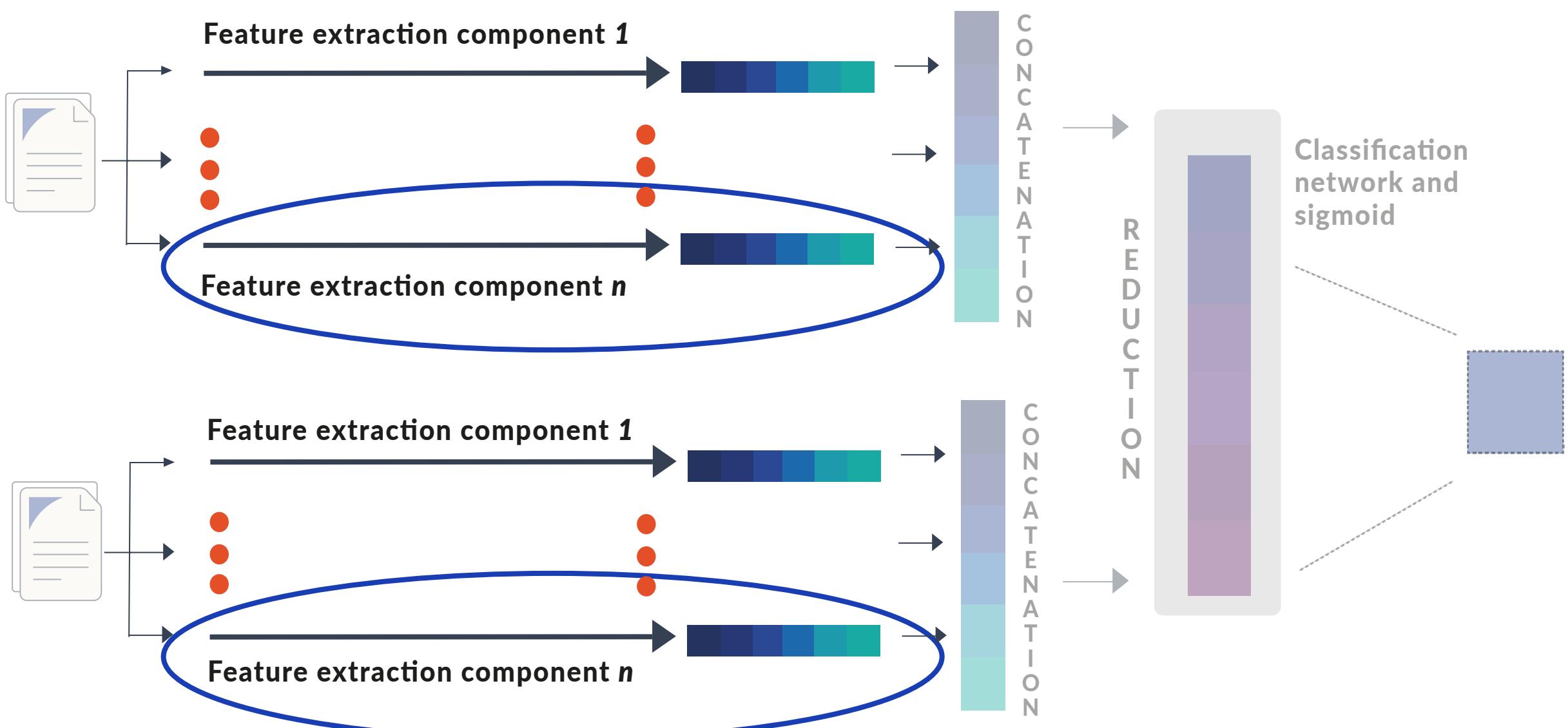
GRAPH-BASED SIAMESE NETWORK ENSEMBLE



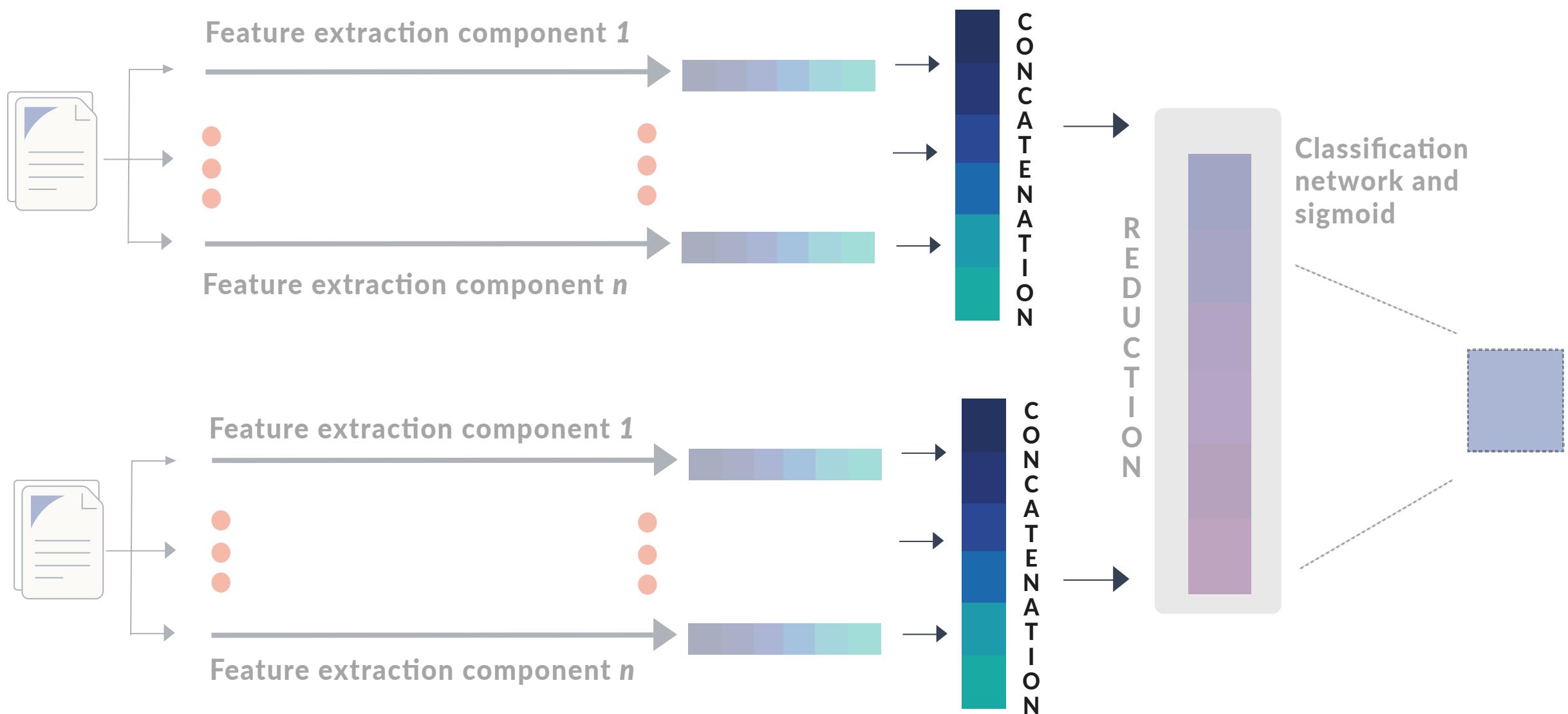
GRAPH-BASED SIAMESE NETWORK ENSEMBLE



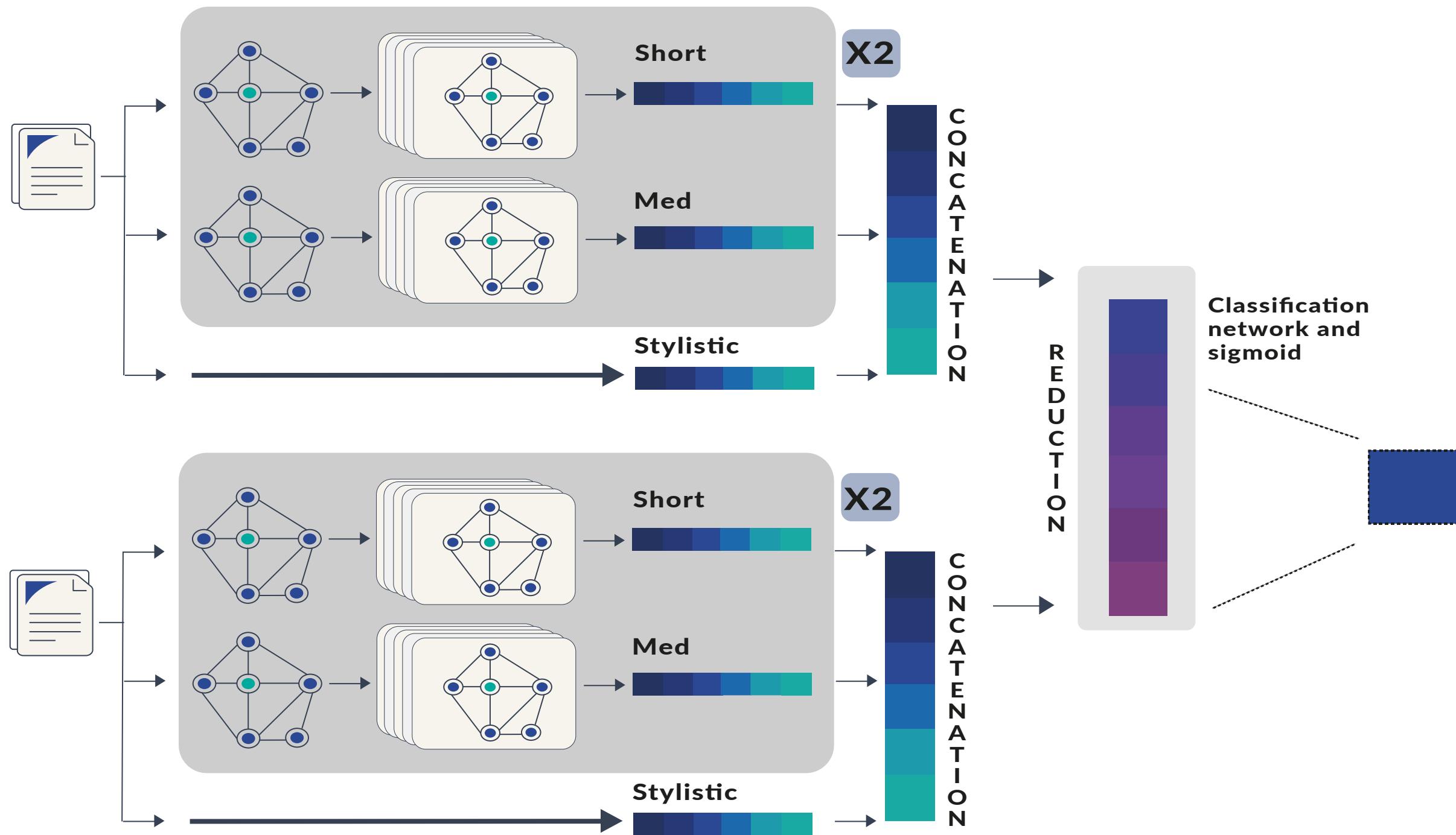
GRAPH-BASED SIAMESE NETWORK ENSEMBLE



GRAPH-BASED SIAMESE NETWORK ENSEMBLE



SUBMITTED MODEL



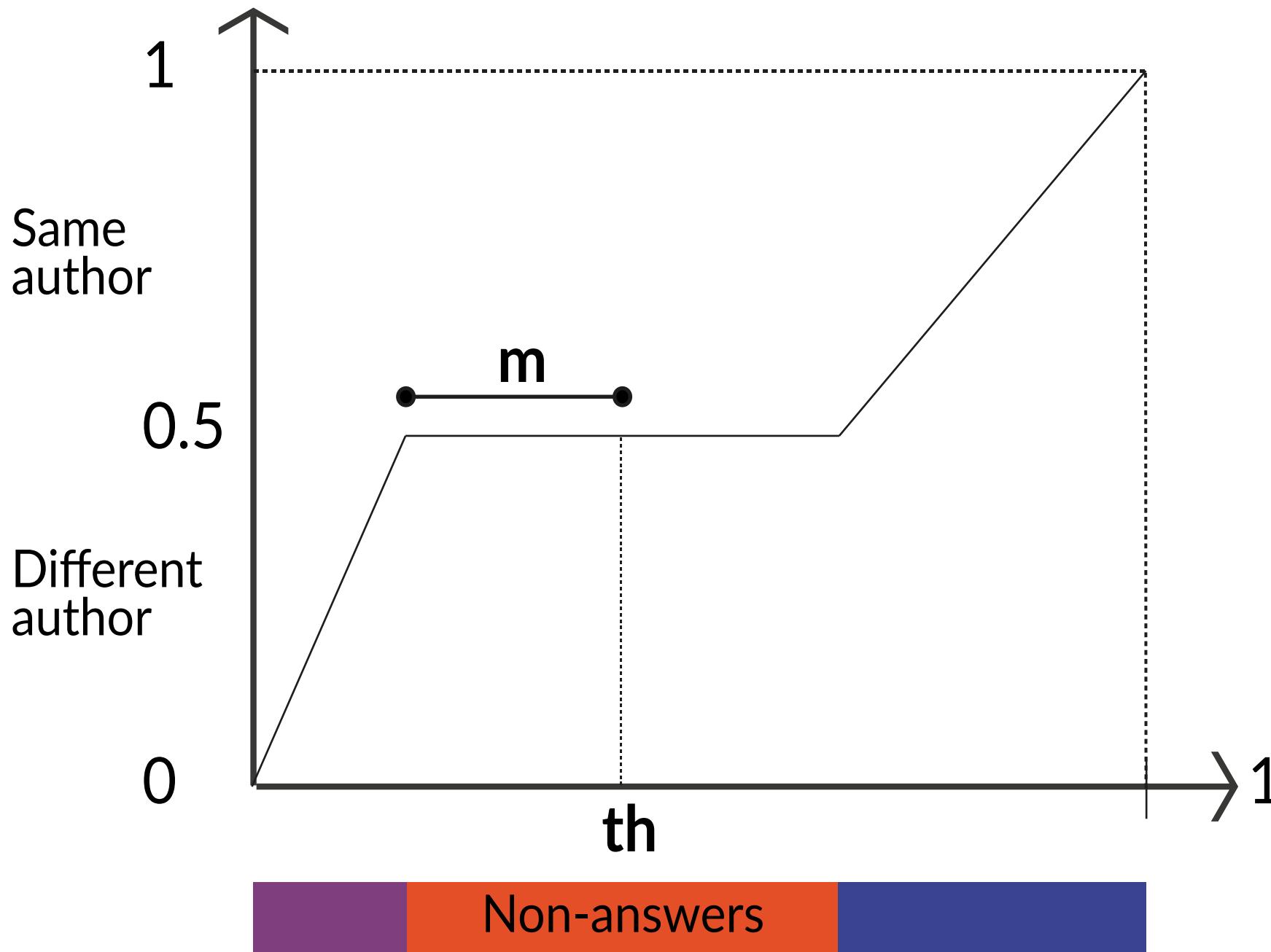
STYLISTIC COMPONENT

- Frequency of **function words**.
- Average number of **characters per word**.
- **Vocabulary richness**.
- Frequency of **words of fixed lenght** (1 to 10 characters).

* Selected from the ones used in:

J Weerasinghe and R. Greenstadt. “Feature Vector Difference based Neural Network and Logistic Regression Models for Authorship Verification”. Notebook for PAN at CLEF 2020, page 6, 2020.

THRESHOLD ADJUSTMENT



th = threshold
 m = margin

We perform a grid search to find the best th and v accordingly to our metrics.

RESULTS

	ARCHITECTURE	SMALL	LARGE
DEVELOPMENT RESULTS	Short graph component	87.81	90.63
	Meg graph component	88.31	91.50
	Short + Med	89.21	92.38
	Short + Med + Features	89.80	NA
	Short(x2) + Med(x2) + Features (Submitted)	90.03	92.96
FINAL RESULT	Short(x2) + Med(x2) + Features (Submitted)	90.70	93.59

CONCLUSIONS

- We propose a **new approach** to the task, using **graph representations** of the text, **graph neural networks** and **siamese networks**.
- Our models were the **2nd** and **7th** best-evaluated models.

FUTURE WORK

- Evaluate other graph representations (REDUCE_LABEL set).
- Evaluate our approach in other tasks.
- Improve the existing model considering character-level information.

THANK YOU!

Daniel Embarcadero Ruiz

danielemburu@gmail.com

Dra Helena Gomez Adorno

helena.gomez@iimas.unam.mx