

Etape 4 : Construction d'un modèle en passant par des résumés de section

La démarche ici est la suivante :

- pour chaque section de brevet on réalise un résumé à base des 5 phrases les plus importantes de la section
- on réalise l'embedding de la section par mean pooling des embeddings des phrases résumant la section
- on réalise l'embedding du brevet par mean pooling des embeddings des résumés de section

```
In [1]: import json
import matplotlib.pyplot as plt
from tqdm import tqdm
import numpy as np
import warnings
warnings.filterwarnings("ignore")
import pickle

from sentence_transformers import SentenceTransformer
from sentence_transformers.util import cos_sim
from sentence_transformers.quantization import quantize_embeddings
from sentence_transformers import losses
from sentence_transformers.readers import InputExample
from torch.utils.data import DataLoader
from transformers import AutoTokenizer

import nltk
import numpy as np
from LexRank import degree centrality_scores
```

Création des embeddings après avoir appliqué des résumés de chacune des sections

Le but ici est d'utiliser un LLM pour créer des résumés de chaque section, puis de créer les embeddings des résumés

```
In [2]: with open('../data/dataset_patent_sections.json', 'r') as outfile:
dataset_patent_section = json.load(outfile)
outfile.close()
```

```
In [21]: model_name = 'intfloat/e5-small-v2'
model = SentenceTransformer(model_name)
```

```
In [84]: # Construction du dataset de brevets resumes par section
```


