

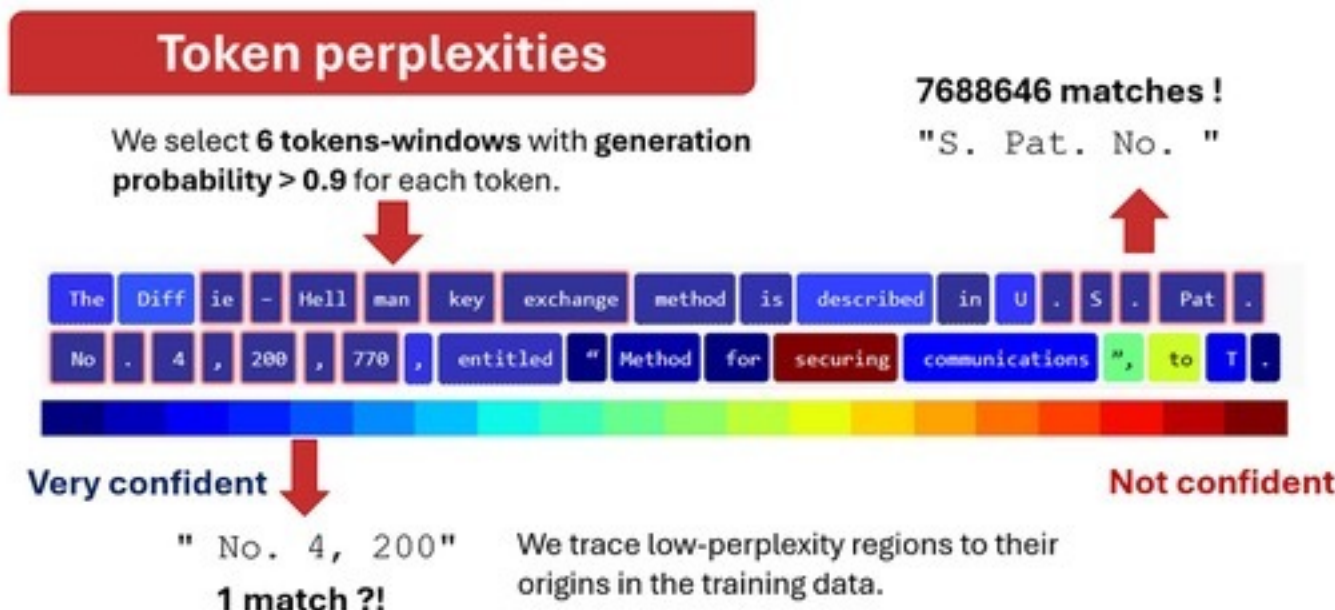


# Low-Perplexity LLM-Generated Sequences and Where To Find Them

Arthur Wuhrmann, Anastasiia Kucherenko, Andrei Kucharavy

**TDA** (Training data attribution) is a field trying to **map sentences to the training data**, particularly with LLMs. This is crucial to copyright checking, data cleaning and to ensure privacy.

When LLMs generate sequences, each token has some probability of being generated. Sometimes, the model is extremely confident in a sequence it generates.



Do LLMs learn passages from their training data **by heart** ?

