

# Introduction to LLM

Practice Session 4

Word Representation I

# Topics Covered Before LLMs

- LLMs implicitly learn classic NLP tasks (POS, parsing, disambiguation)
  - We teach them first to show what LLMs do implicitly, e.g., morphology, handling inflection, OOV.
  - These topics show what's happening under the hood.
- You learn how text becomes numbers and how we search and compare with it, key ideas behind modern LLM systems.
- Seeing pre-LLM methods teaches trade-offs and when simple tools are enough.
- Not every problem needs an LLM, always start with a baseline.
  - For small, well-defined, labeled task LLMs will be an overkill (e.g., 2-layer NN for tweet sentiment can hit  $\geq 90\%$ ).
- Evaluate fairly (simple models vs. LLM) to justify cost, latency, and privacy impacts.
- Takeaway: Understand the classics, then choose the right tool based on task complexity.
  - Use LLMs only when gains justify cost or complexity.

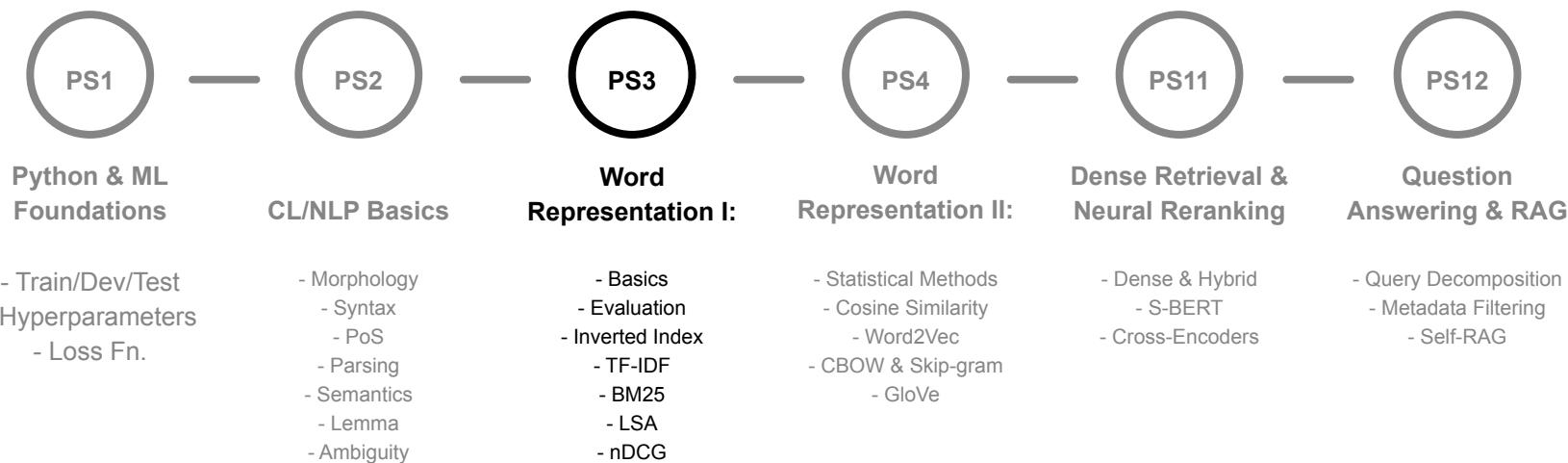
# Clarifying Lesk Variants

- Original Lesk (1986)
  - Algorithm explained in Lecture
  - Compare glosses of all candidate senses with glosses of surrounding words' senses; score by overlap between definitions.
  - *(definitions ↔ definitions overlap)*
- Simplified Lesk (Kilgarriff & Rosenzweig (2000)):
  - Algorithm explained in PS & HW
  - Compare sentence context directly to each sense's gloss.
  - *(definitions ↔ sentence context overlap)*
- Simplified Lesk is the one you'll see used most often in practice (as a baseline and in variants):
  - Mainly because it's far cheaper computationally and tends to work better than the original on standard evaluations.

# Difficulty of Exercise and HW

- It's normal to feel a bit overwhelmed in PS: moving from theory to code isn't straightforward.
- There's a natural gap: you may understand the lecture but struggle in implementation
  - PS is designed to bridge that gap.
- This tries to mirror your first job: turning ideas/specs into working, efficient code for your problem.
- PS focus: review the concepts in lecture and see how can we implement a minimal working version.
  - After that we show you how you can use reliable libraries to avoid reinventing the wheel.
- Important that you have programming basics + ML foundations (we tried to review that in PS1).
  - We recommended external resources on Moodle, e.g. 3Blue1Brown YT channel for visual intuition.
- It's okay if you don't understand everything live.
  - Ideally you understand the main idea and purpose of the code, then you can go over the details step by step yourself.

# Timeline



# PS3: Colab Notebook (Available on Moodle)



The screenshot shows a Google Colab interface with the following details:

- Title:** Copy of Practice\_Session\_01\_Student\_Exercises.ipynb
- Toolbar:** File, Edit, View, Insert, Runtime, Tools, Help
- Search Bar:** Commands
- Code Editor:** A large text area containing the notebook content.
- Toolbar Buttons:** + Code, + Text, Run all, Copy to Drive (highlighted with a dashed red box and arrow).
- Table of Contents:**
  - PS 01: Introduction (Python and ML Foundations)
  - Learning Objectives
- Description:** By the end of this practice session, you will be able to:
- Objectives List:**
  1. Know Python basics: variables, data types, operators, and control structures
  2. Manipulate strings effectively: indexing, slicing, and built-in string methods
  3. Work with data structures: lists, dictionaries, and their operations
  4. Handle file I/O: reading/writing text files and JSON data
  5. Use essential libraries: NumPy for numerical computing, Pandas for data manipulation
  6. Apply object-oriented programming: classes, methods, and type hints
  7. Implement basic ML workflows: data splitting, model training with PyTorch

- <https://colab.research.google.com/drive/1IdxJCu6HOsoOJUXwnCJ6-jrF-ypjGL9G>