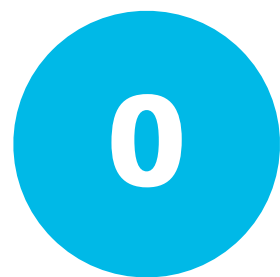Natural Language Processing
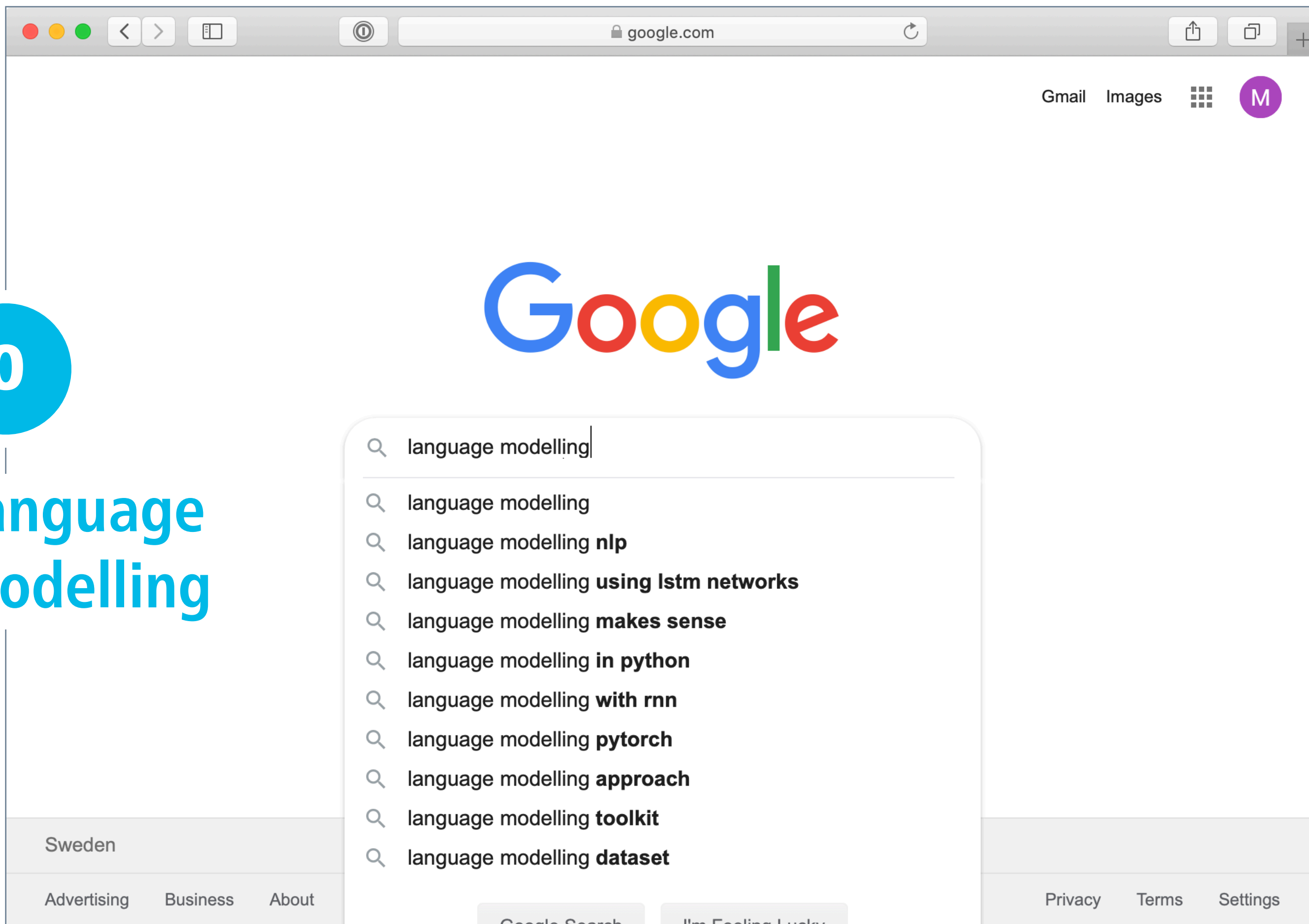
# Course overview

Marco Kuhlmann

Department of Computer and Information Science
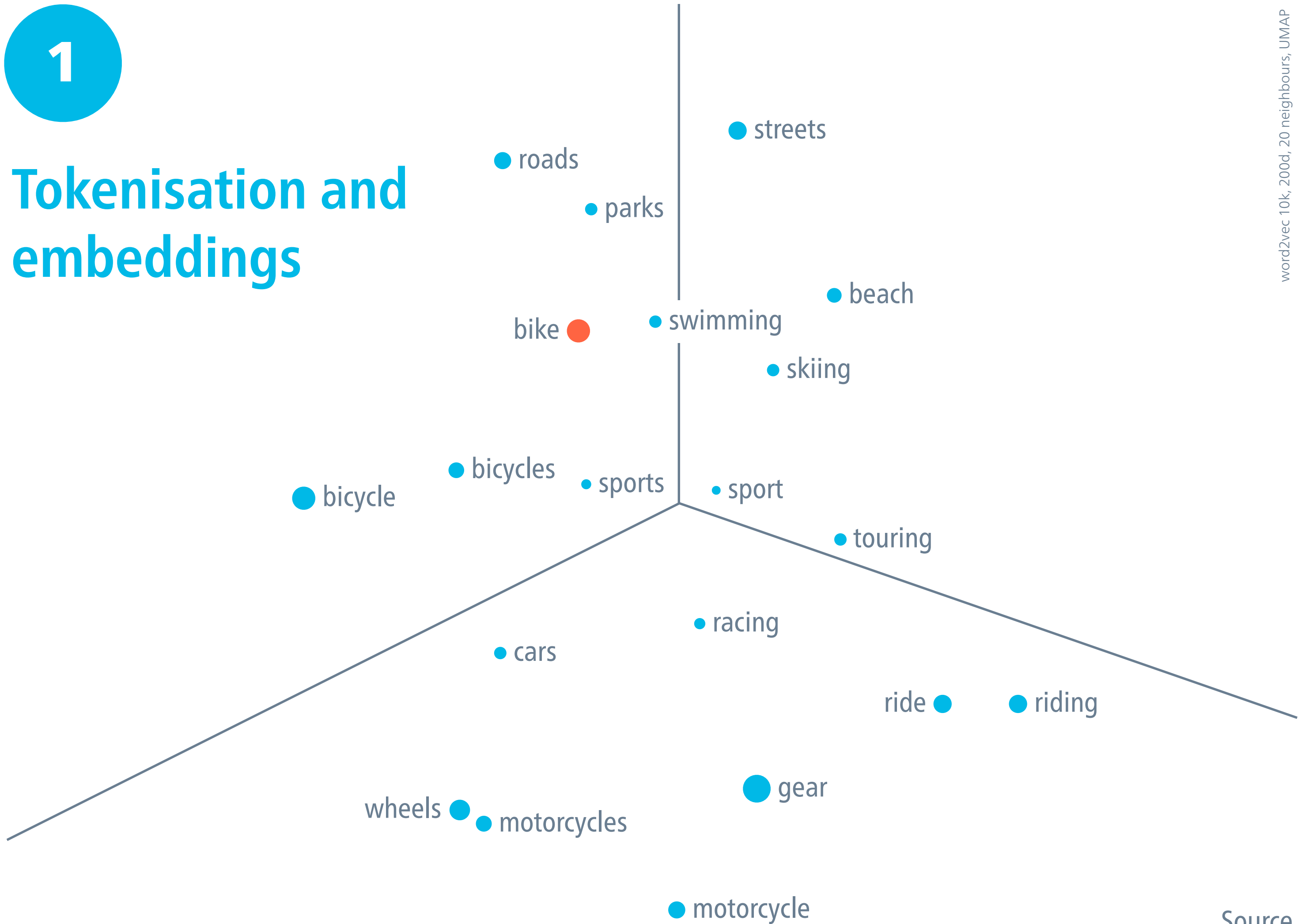
**0**

**Language modelling**

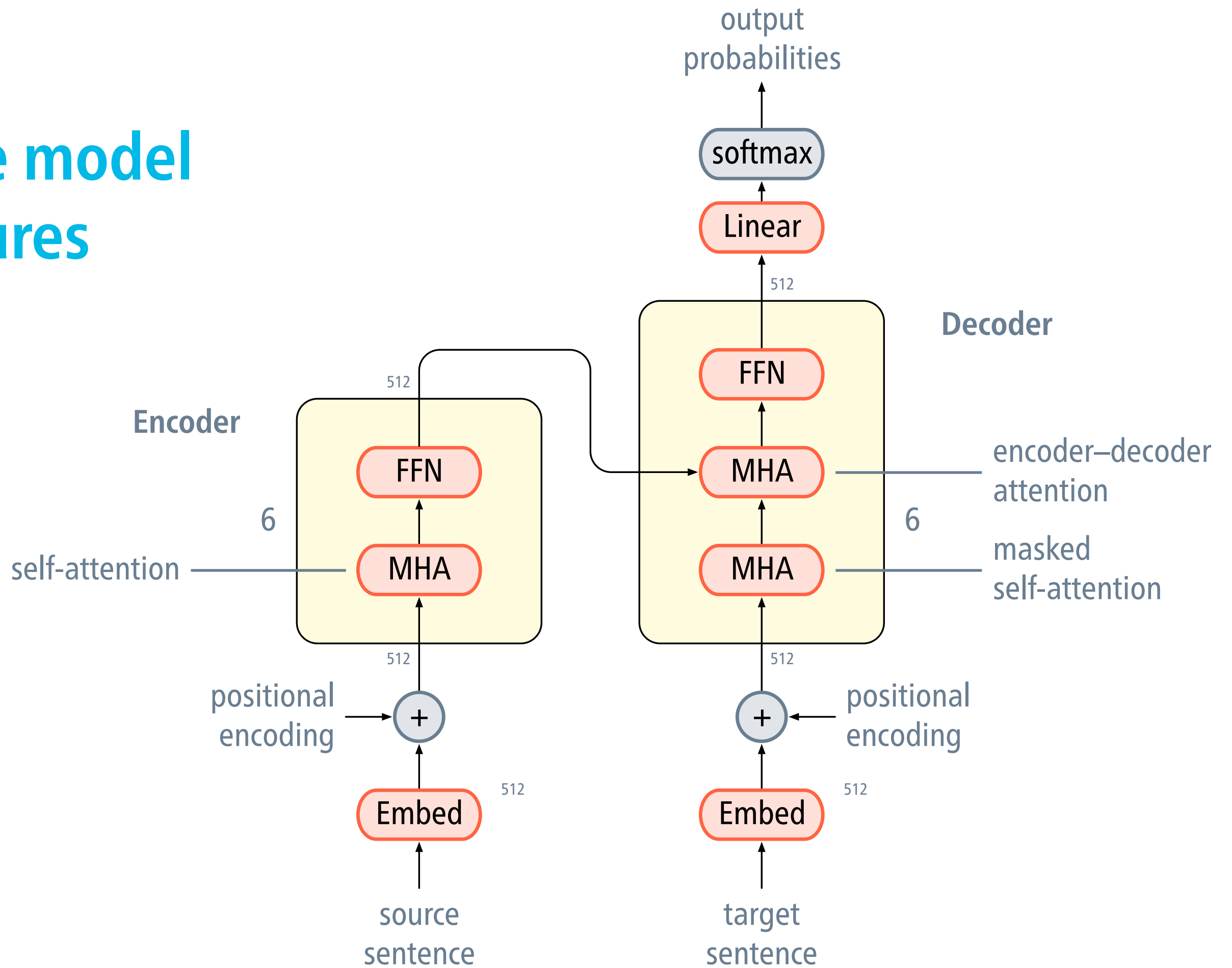# Tokenisation and embeddings

- streets
- roads
- parks
- beach
- bike
- swimming
- skiing
- bicycles
- bicycle
- sports
- sport
- touring
- racing
- cars
- ride
- riding
- gear
- wheels
- motorcycles
- motorcycle

Source

**2**

# Language model architectures

output probabilities

source sentence

target sentence

Encoder

Decoder

self-attention

encoder–decoder attention

masked self-attention

softmax

Linear

FFN

MHA

FFN

MHA

MHA

Embed

Embed

positional encoding

positional encoding

+

+

512

512

512

512

512

512

512

6

6

**3** Pre-training and fine-tuning

# 5 Current research

## Paper 1 (left)

**Ignore This Title and HackAPrompt: Exposing Systemic Vulnerabilities of LLMs through a Global Scale Prompt Hacking Competition**

Sander Schulhoff[1*]  Jeremy Pinto[2*]  Anaum Khan[1]  Louis-François Bouchard[2,3]  Chr...
Svetlina Anati[5**]  Valen Tagliabue[6**]  Anson Liu Kost[7**]  Christopher Carnahan...
Jordan Boyd-Graber[1]
[1] University of Maryland  [2] Mila  [3] Towards AI  [4] Stanford
[5] Technical University of Sofia  [6] University of Milan  [7] NYU
[8] University of Arizona
sschulho@umd.edu  jerpint@gmail.com  jbg@umiacs.umd.edu

**Abstract**

Large Language Models (LLMs) are deployed in interactive contexts with direct user engagement, such as chatbots and writing assistants. These deployments are vulnerable to prompt injection and jailbreaking (collectively, prompt hacking), in which models are manipulated to ignore their original instructions and follow potentially malicious ones. Although widely acknowledged as a significant security threat, there is a dearth of large-scale resources and quantitative studies on prompt hacking. To address this lacuna, we launch a global prompt hacking competition, which allows for free-form human input attacks. We elicit 600K+ adversarial prompts against three state-of-the-art LLMs. We describe the dataset, which empirically verifies that current LLMs can indeed be manipulated via prompt hacking. We also present a comprehensive taxonomical ontology of the types of adversarial prompts.

### 1 Introduction: Prompted LLMs are Everywhere... How Secure are They?

Large language models (LLMs) instruct these LLM models what to do (Zamfirescu-Pereira et al., 2023; Khashabi et al., 2022; Min et al., 2022; Webson and Pavlick, 2022). The flexibility of this approach not
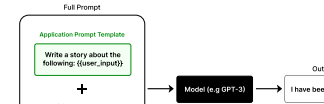
Figure 1: Uses of LLMs often define the task in a prompt template (top left), which is combined with input (bottom left). We create a competition to...

4945

## Paper 2 (center)

**Label Words are Anchors: An Information Flow Perspective for Understanding In-Context Learning**

Lean Wang[†,§], Lei Li[†], Damai Dai[†], Deli Chen[§],
Hao Zhou[§], Fandong Meng[§], Jie Zhou[§], Xu Sun[†]
[†]National Key Laboratory for Multimedia Information Processing,
School of Computer Science, Peking University
[§]Pattern Recognition Center, WeChat AI, Tencent Inc., China
{lean,daidai,xusun}@pku.edu.cn  nlp.lilei@gmail.com
victorchen@deepseek.com  {tuxzhou,fandongmeng,withtomzhou}@tencent.com

**Abstract**

In-context learning (ICL) emerges as a promising capability of large language models (LLMs) by providing them with demonstration examples to perform diverse tasks. However, the underlying mechanism of how LLMs learn from the provided context remains under-explored. In this paper, we investigate the working mechanism of ICL through an information flow lens. Our findings reveal that label words in the demonstration examples function as anchors: (1) semantic information aggregates into label word representations during the shallow computation layers' processing; (2) the consolidated information in label words serves as a reference for LLMs' final predictions. Based on these insights, we introduce an anchor re-weighting method to improve ICL performance, a demonstration compression technique to expedite inference, and an analysis framework for diagnosing ICL errors in GPT2-XL. The promising applications of our findings again validate the uncovered ICL working mechanism and pave the way for future studies.

Figure 1: Visualization of the information flow in a GPT model performing ICL. The line depth reflects the significance of the information flow from the right word to the left. The flows involving label words are highlighted. Label words gather information from demonstrations in shallow layers, which is then extracted in deep layers for final prediction.[1]

### 1 Introduction

In-Context Learning (ICL) has emerged as a powerful capability alongside the development of scaled-up large language models (LLMs) (Brown et al., 2020). By instructing LLMs using few-shot demonstration examples, ICL enables them to perform a wide range of tasks, such as text classification (Min et al., 2022a) and mathematical reasoning (Wei et al., 2022). Since ICL does not require updates to millions or trillions of model parameters and relies on human-understandable natural language instructions (Dong et al., 2023), it has become a promising approach for harnessing the full potentiality of LLMs. Despite its significance, the inner working mechanism of ICL remains an open question, garnering considerable interest from research communities (Xie et al., 2022; Dai et al., 2022; Akyürek et al., 2022; Li et al., 2023b).

In this paper, we find that the label words serve as anchors that aggregate and distribute information in ICL. We first visualize the attention interactive pattern between tokens with a GPT model (Brown et al., 2020) on sentiment analysis (Figure 1). Initial observations suggest that label words aggregate information in shallow layers and distribute it in deep layers.[2] To draw a clearer picture of this phenomenon, we design two metrics based on saliency

[1] https://github.com/lancopku/label-words-are-anchors

[2] In this paper, "shallow" or "first" layers refer to those closer to the input, while "deep" or "last" layers are closer to the output. Here, "deep layers" include those around the midpoint, e.g., layers 25-48 in a 48-layer GPT2-XL.

9840

## Paper 3 (right)

**...ster Minimum Bayes Risk Decoding with Confidence-based Pruning**

Julius Cheng, Andreas Vlachos
Department of Computer Science and Technology
University of Cambridge
{jncc3,av308}@cam.ac.uk

**Abstract**

Minimum Bayes risk (MBR) decoding outputs the hypothesis with the highest expected utility over the model distribution for some utility function. It has been shown to improve accuracy over beam search in conditional language generation problems and especially neural machine translation, in both human and automatic evaluations. However, the standard sampling-based algorithm for MBR is substantially more computationally expensive than beam search, requiring a large number of samples as well as a quadratic number of calls to the utility function, limiting its applicability. We describe an algorithm for MBR which gradually grows the number of samples used to estimate the utility while pruning hypotheses that are unlikely to have the highest utility according to confidence estimates obtained with bootstrap sampling. Our method requires fewer samples and drastically reduces the number of calls to the utility function compared to standard MBR while being statistically indistinguishable in terms of accuracy. We demonstrate the effectiveness of our approach in experiments on three language pairs, using chrF++ and COMET as utility/evaluation metrics.

### 1 Introduction

Minimum Bayes risk (MBR) decoding (Bickel and Doksum, 1977; Goel and Byrne, 2000) has recently received renewed attention as a decision rule for conditional sequence generation tasks, especially machine translation (NMT). In MBR, the sequence with the highest expected utility with respect to the model distribution is chosen as the output, where the utility is usually some measure of similarity. This contrasts with the more commonly used maximum a posteriori (MAP) decision rule, which returns the sequence with the highest probability under the model. MAP is generally intractable, and beam search is typically used to find an approximation. MBR is likewise intractable,

### 2 Minimum Bayes risk decoding

Conditional sequence generation problems such as neural machine translation (NMT) model the probability of the next token $y_t$ given a source sequence $x$ and prefix $y_{<t}$ with a neural network $p_\theta$. This

12473