**Faculty of Engineering, Built Environment and Information Technology**

Fakulteit Ingenieurswese, Bou-omgewing en Inligtingtegnologie / Lefapha la Boetšenere, Tikologo ya Kago le Theknolotši ya Tshedimošo
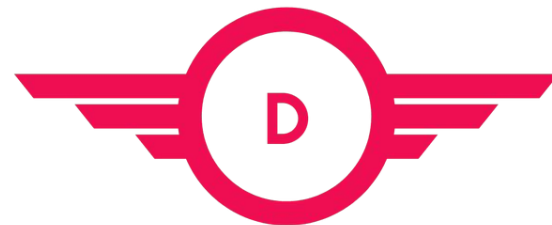
UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

# Special Topic: Word Embeddings + Language Models

Dr. Vukosi Marivate

**Inputs:**
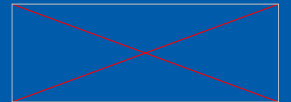- A Modupe [PhD Candidate]
- A Moodley [MIT Big Data Science Student]

Make today matter

Data Science for Social Impact

# Word Embeddings + Language Models

# Word Embeddings

**"You shall know a word by the company it keeps"**

**"Tell me who your friends are, and I will tell you who you are."**

**The Distributional Hypothesis** is that words that occur in the same contexts tend to have similar meanings [2]
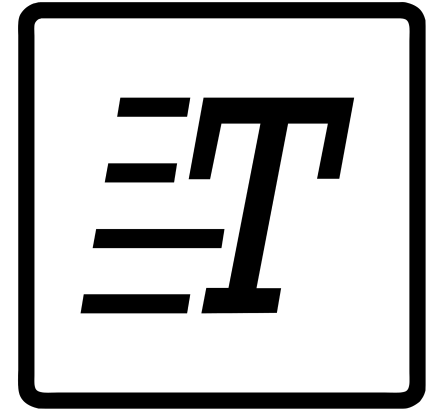
# Word Embeddings

**"You shall know a word by the company it keeps"**

**"Tell me who your friends are, and I will tell you who you are."**

**The Distributional Hypothesis** is that words that occur in the same contexts tend to have similar meanings [2]
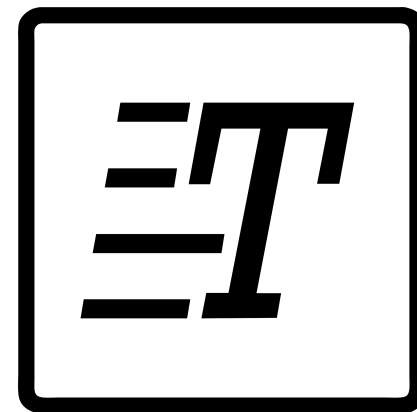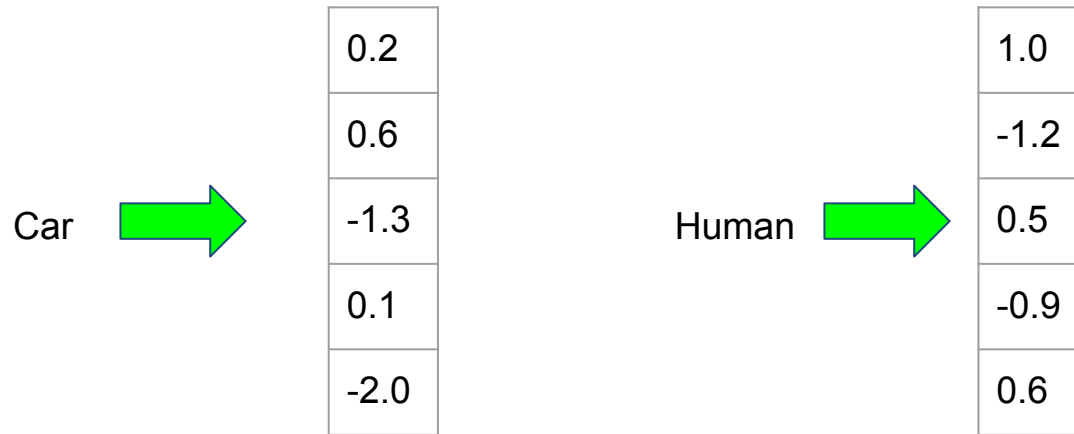
# Word Vectors

- Mapping from tokens to a continuous vector space

- Trained using a shallow neural network (not deep)

Car →

| 0.2 |
|------|
| 0.6 |
| -1.3 |
| 0.1 |
| -2.0 |

Human →

| 1.0 |
|------|
| -1.2 |
| 0.5 |
| -0.9 |
| 0.6 |

UNIVERSITEIT VAN PRETORIA
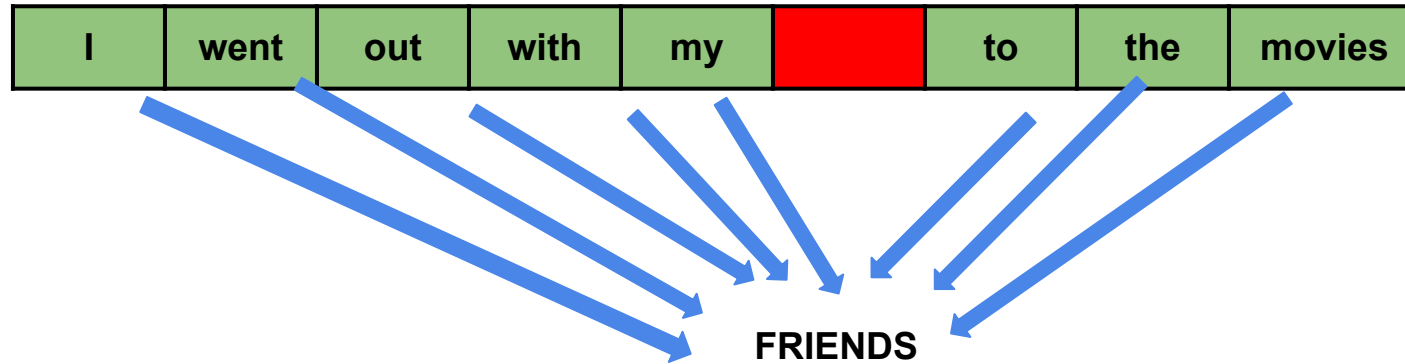UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

# Word Vectors - The Idea!

- I went out with my **BOYFRIEND** to the movies.
- I went out with my **GIRLFRIEND** to the movies.
- I went out with my **BAE** to the movies.
- I went out with my **FRIENDS** to the movies.

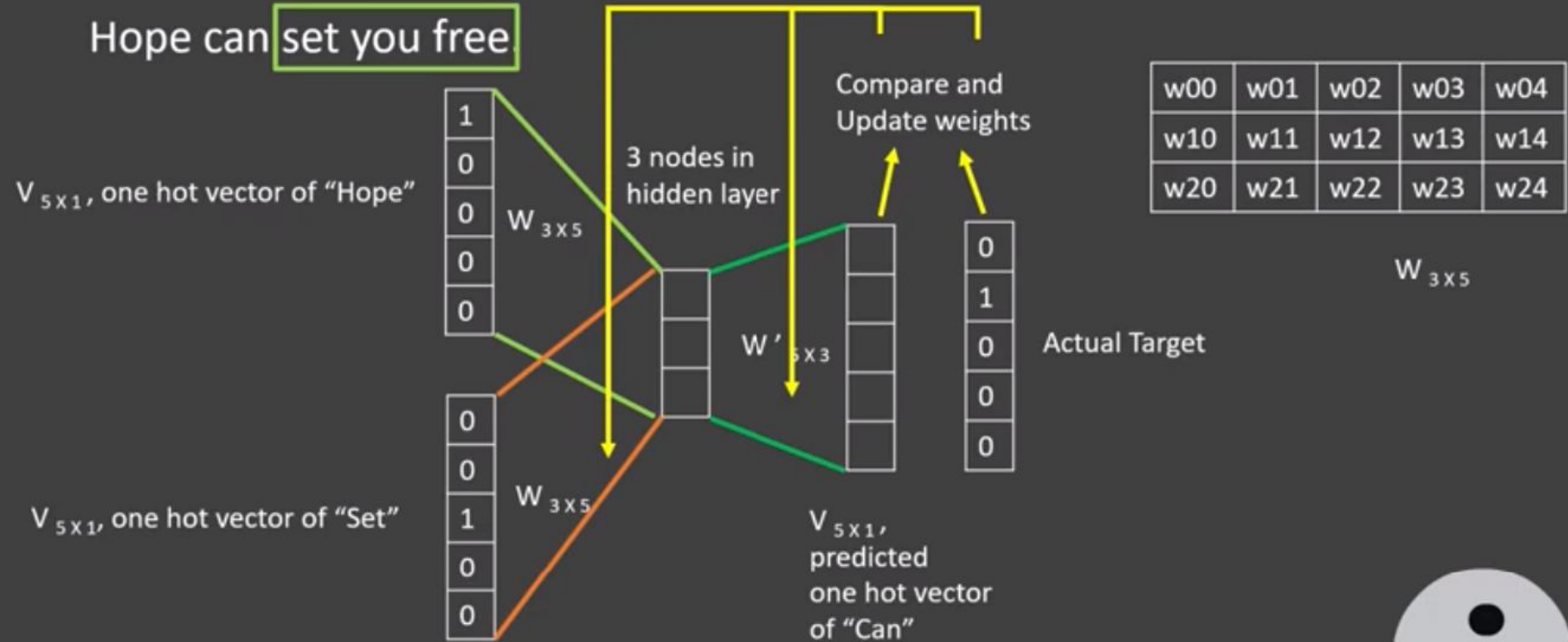Words used in a similar fashion in the same context!!!!

https://www.youtube.com/watch?v=UqRCEmrv1gQ

UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

# CBOW - Continuous Bag of Words

Predict word from context



| I | went | out | with | my | | to | the | movies |
|---|------|-----|------|----|----|----|----|--------|

**FRIENDS**

UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

# CBOW - Working

Hope can set you free.

$V_{5 \times 1}$, one hot vector of "Hope"

$\begin{array}{|c|}\hline 1 \\\hline 0 \\\hline 0 \\\hline 0 \\\hline 0 \\\hline\end{array}$   $W_{3 \times 5}$

3 nodes in hidden layer

$W'_{5 \times 3}$

$V_{5 \times 1}$, one hot vector of "Set"

$\begin{array}{|c|}\hline 0 \\\hline 0 \\\hline 1 \\\hline 0 \\\hline 0 \\\hline\end{array}$   $W_{3 \times 5}$

Compare and Update weights

$\begin{array}{|c|}\hline \phantom{0} \\\hline \phantom{0} \\\hline \phantom{0} \\\hline \phantom{0} \\\hline \phantom{0} \\\hline\end{array}$

$\begin{array}{|c|}\hline 0 \\\hline 1 \\\hline 0 \\\hline 0 \\\hline 0 \\\hline\end{array}$   Actual Target

$V_{5 \times 1}$, predicted one hot vector of "Can"

| w00 | w01 | w02 | w03 | w04 |
|-----|-----|-----|-----|-----|
| w10 | w11 | w12 | w13 | w14 |
| w20 | w21 | w22 | w23 | w24 |

$W_{3 \times 5}$

YUNIBESITHI YA PRETORIA

# Skip-GRAM: Crazy idea, but works!!

Predict context from words



FRIENDS

| I | went | out | with | my | | to | the | movies |

UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
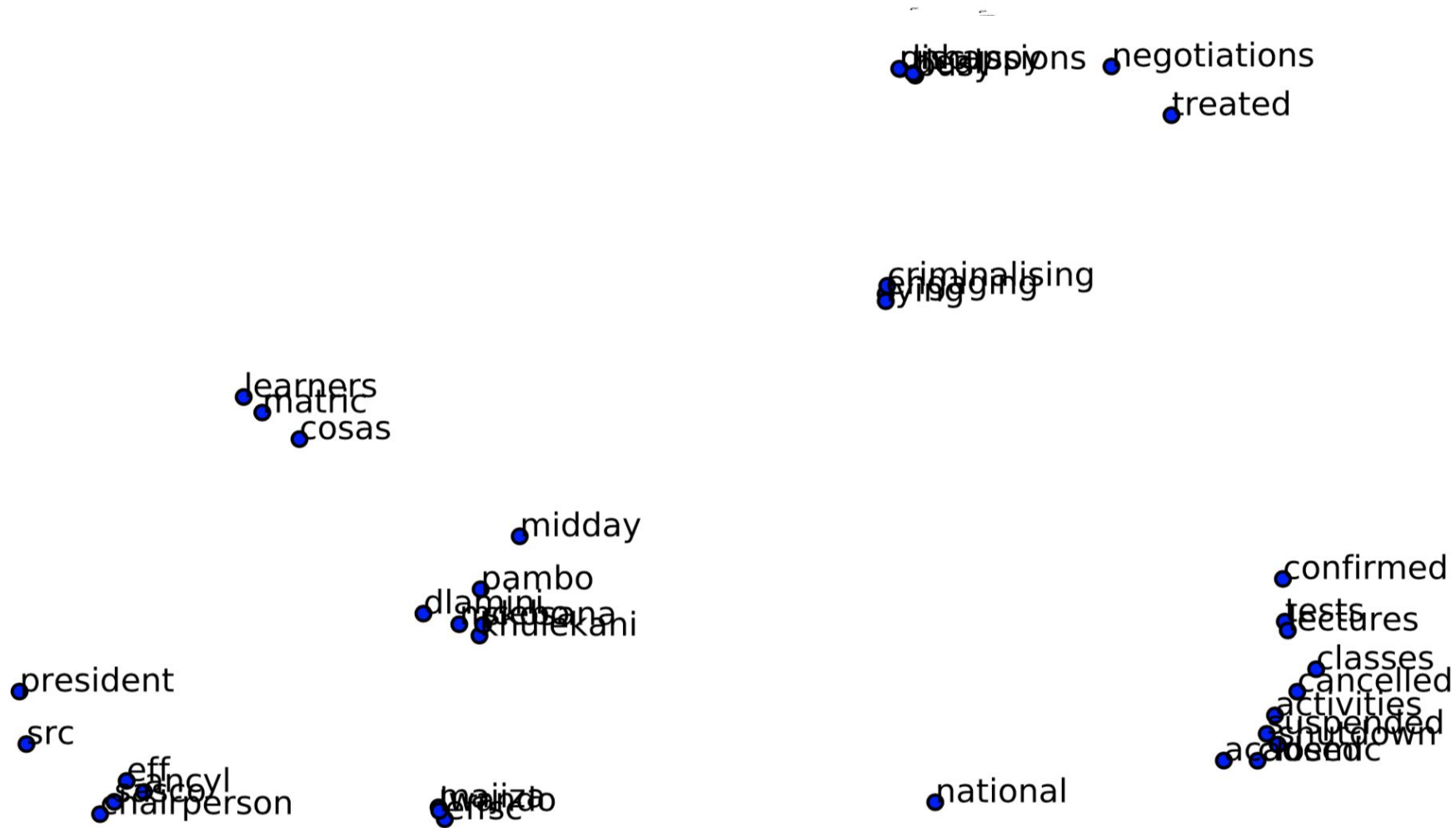YUNIBESITHI YA PRETORIA

# Word2Vec



$$\frac{1}{T}\sum_{t=1}^{T}\sum_{-c\leq j\leq c,j\neq 0} logp(w_t|w_{t+j})$$

$$\frac{1}{T}\sum_{t=1}^{T}\sum_{-c\leq j\leq c,j\neq 0} logp(w_{t+j}|w_t)$$

# GloVe: Global Vectors for Word Representation

**Nearest neighbors**

The Euclidean distance (or cosine similarity) between two word vectors provides an effective method for measuring the linguistic or semantic similarity of the corresponding words. Sometimes, the nearest neighbors according to this metric reveal rare but relevant words that lie outside an average human's vocabulary. For example, here are the closest words to the target word *frog*:

0. *frog*
1. frogs
2. toad
3. litoria
4. leptodactylidae
5. rana
6. lizard
7. eleutherodactylus

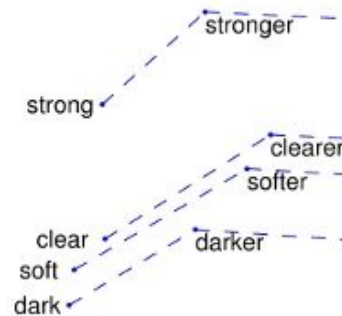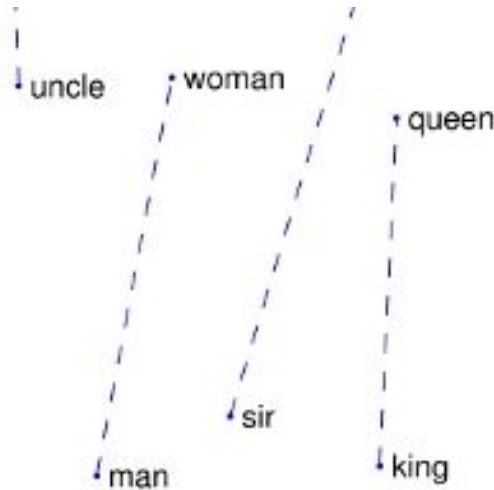3. litoria    4. leptodactylidae    5. rana    7. eleutherodactylus
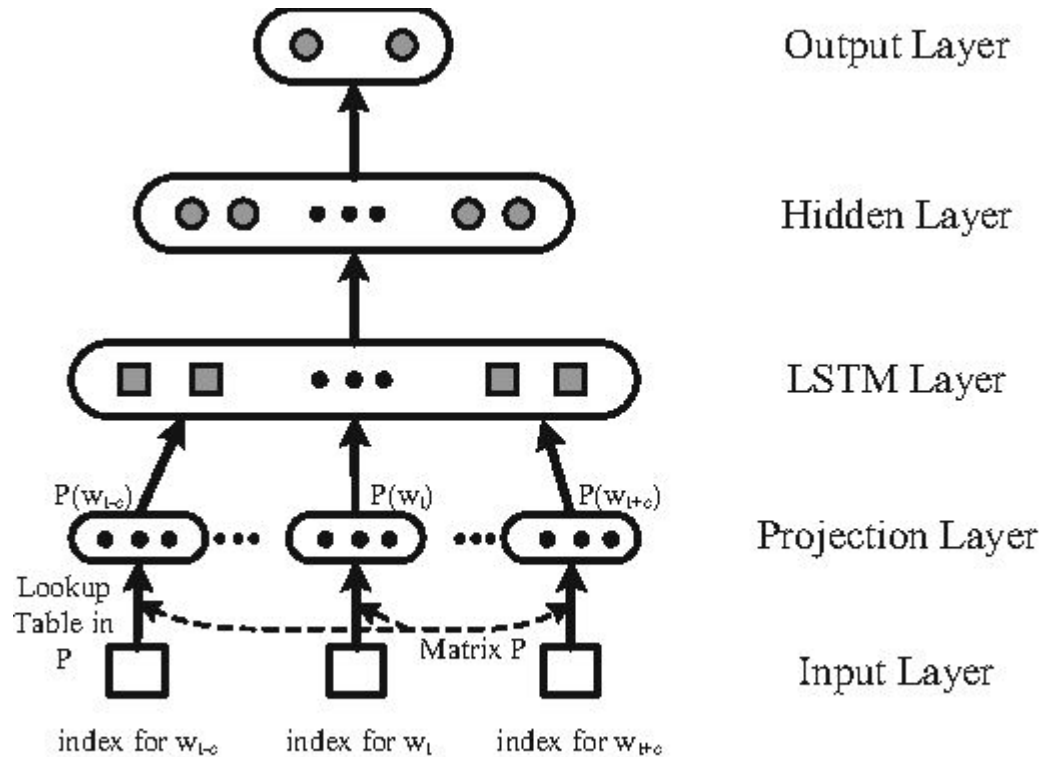
Linear Substructures: Analogies

[Pre built vectors]

UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

13

# Downstream tasks: Classification



Output Layer

Hidden Layer

LSTM Layer

$P(w_{t-c})$    $P(w_t)$    $P(w_{t+c})$    Projection Layer

Lookup Table in P    Matrix P    Input Layer

index for $w_{t-c}$    index for $w_t$    index for $w_{t+c}$

A Bidirectional LSTM Approach with Word Embeddings for Sentence Boundary Detection
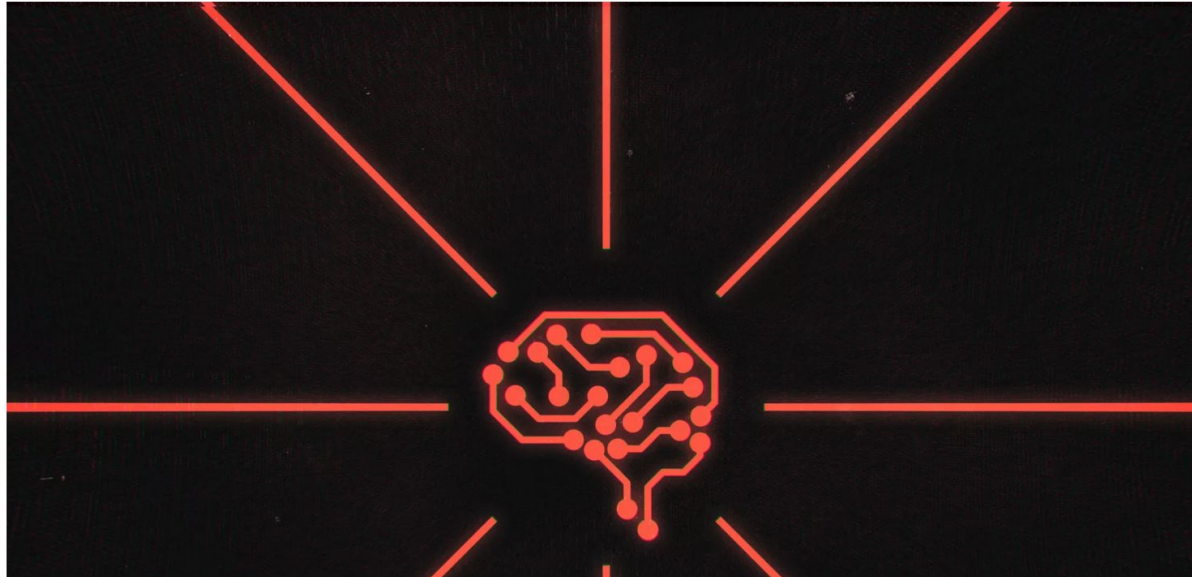
# Advanced: Language Models

## AI researchers debate the ethics of sharing potentially harmful programs

*Nonprofit lab OpenAI withheld its latest research, but was criticized by others in the field*
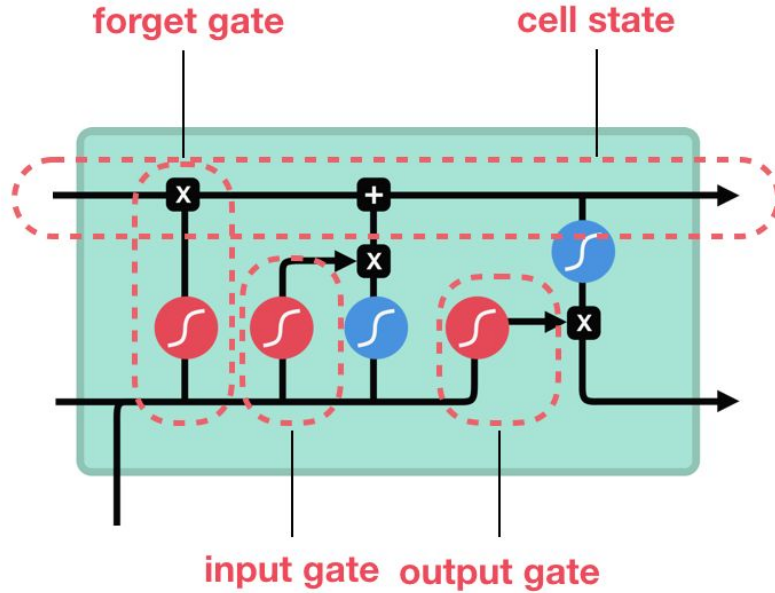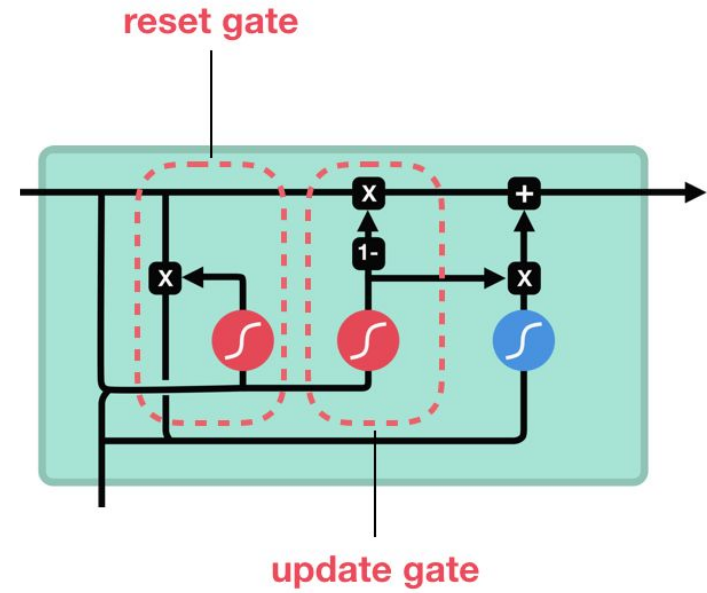
By James Vincent | Feb 21, 2019, 10:30am EST

f ⬤ SHARE

**LSTM**          **GRU**

forget gate     cell state      reset gate

input gate   output gate      update gate

sigmoid     tanh     pointwise multiplication     pointwise addition     vector concatenation

**Illustrated Guide to LSTM's and GRU's: A step by step explanation**

# Sequence Models



**Neural Machine Translation**
SEQUENCE TO SEQUENCE MODEL

Encoding Stage

Decoding Stage

Encoder RNN

Decoder RNN

Je          suis          étudiant

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT   Je    suis   étudiant

UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

Decoding time step: 1 (2) 3 4 5 6    OUTPUT    I

Kencdec    Vencdec

Linear + Softmax

ENCODERS    DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant    PREVIOUS OUTPUTS    I

UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA
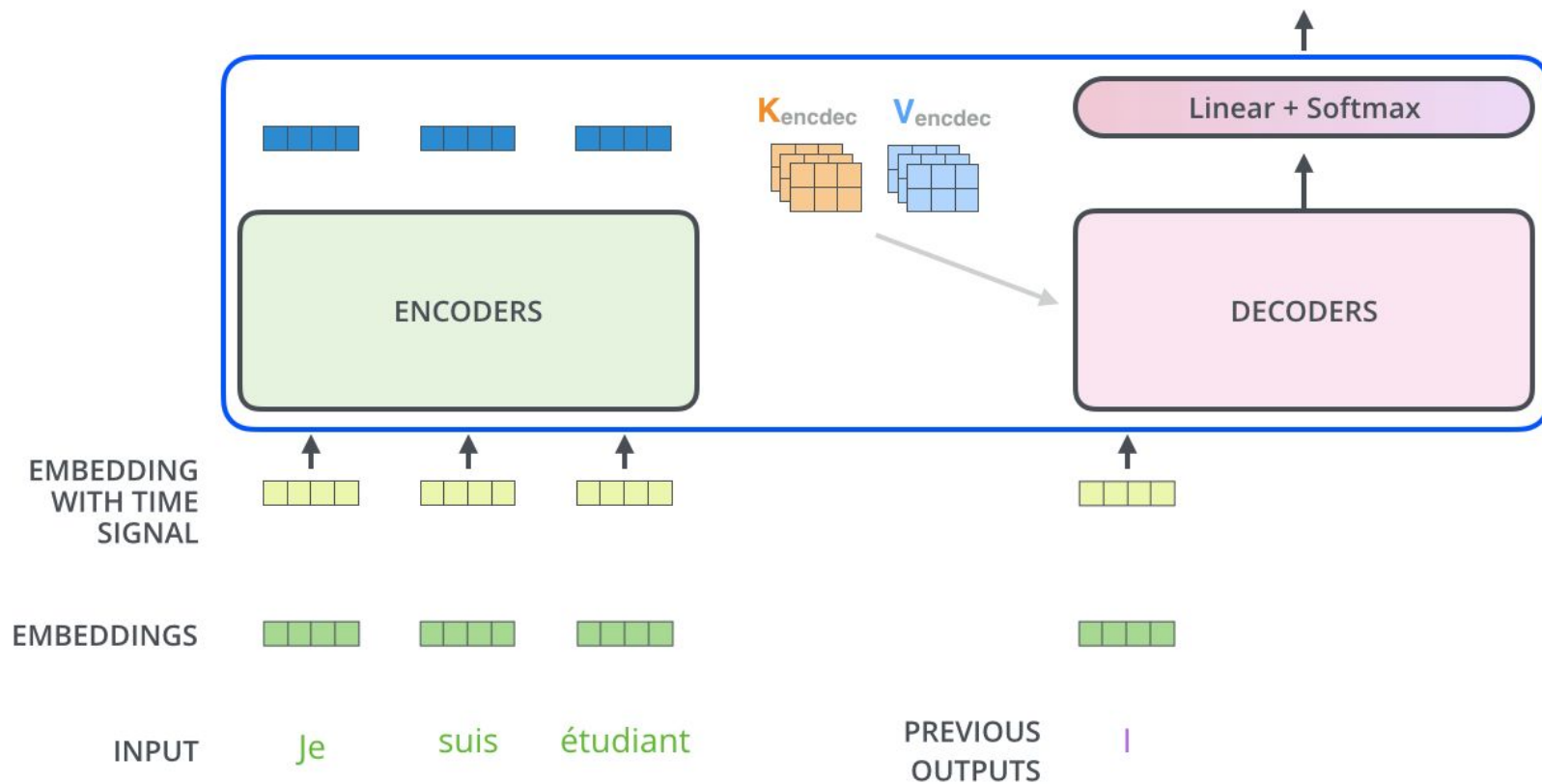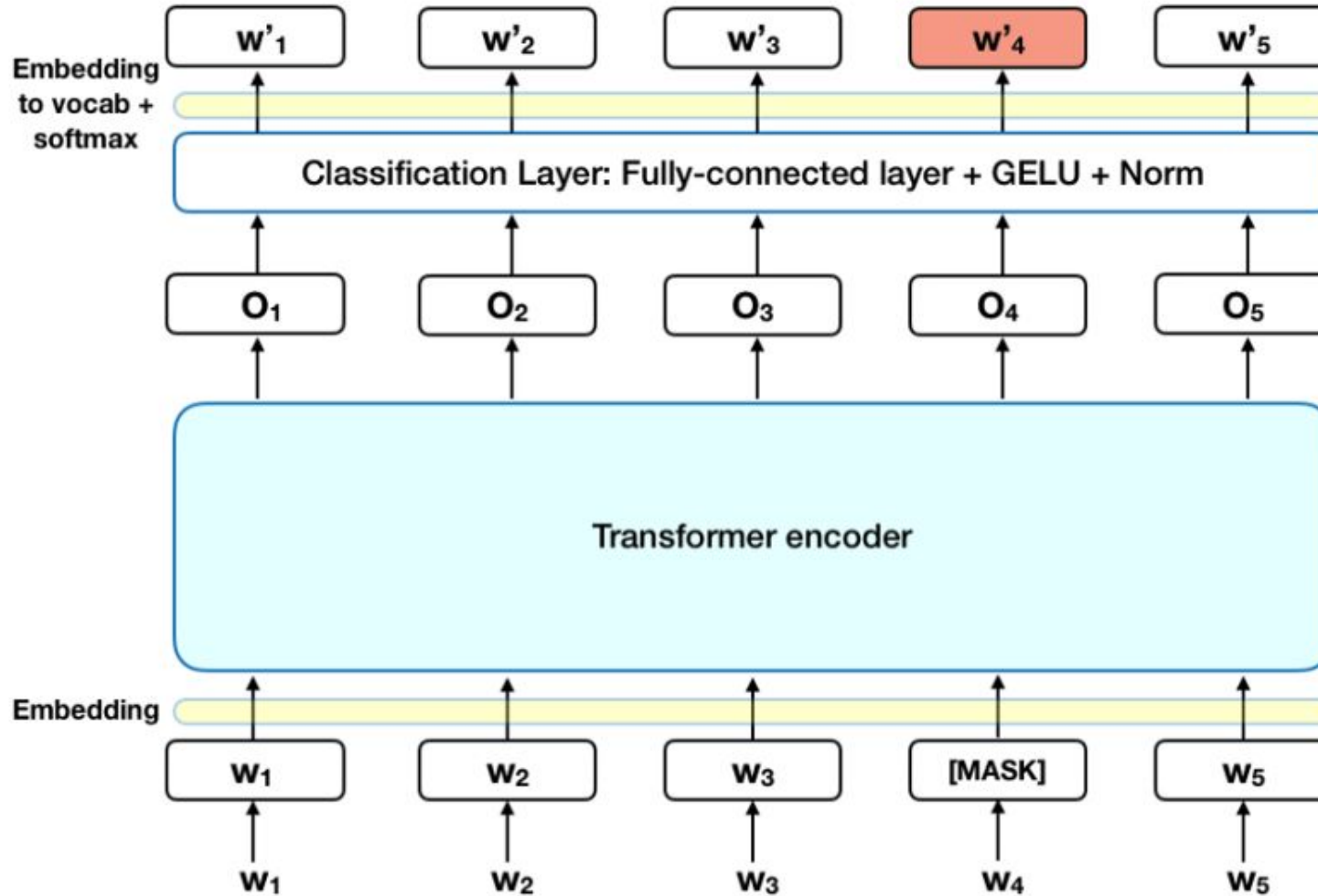
# BERT Masked LM

# What can you do with these LMs

Part of many ML pipelines

Natural Language Understanding

Question Answering

Text Generation [https://transformer.huggingface.co/]

Topic Models

# Transf



Domain adaptation

Transductive transfer learning

Different domains

Same task; labeled data

**Pretraining**

word2vec
GloVe
skip-thought
InferSent
ELMo
ULMFiT
GPT
BERT

**Adaptation**

classification
sequence labeling
Q&A
....

Tasks learned sequentially

Sequential transfer learning

A taxonomy for transfer learning in NLP (Ruder, 2019).

UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

# Resources

**NLP General**

- https://github.com/fastai/course-nl
- Stanford Coursera NLP Slides
  https://web.stanford.edu/~jurafsky/NLPCourseraSlides.html
- Sebastian Ruder Newsletter [http://ruder.io/nlp-news/]
- https://nlpprogress.com

**Python Libraries**

- SKLearn NLP (Working With Text Data) - URL (Nice tutorial)
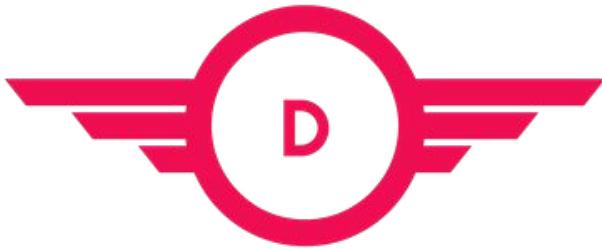- spaCY: Industrial-Strength Natural Language Processing - URL
- NLTK

# Thank You

**Dr. Vukosi Marivate**

vukosi.marivate@cs.up.ac.za

https//dsfsi.github.io

@vukosi



Data Science for Social Impact