# Probing as a Technique to Understand Abstract Spaces

Ashlen Plasek
University of Minnesota Morris
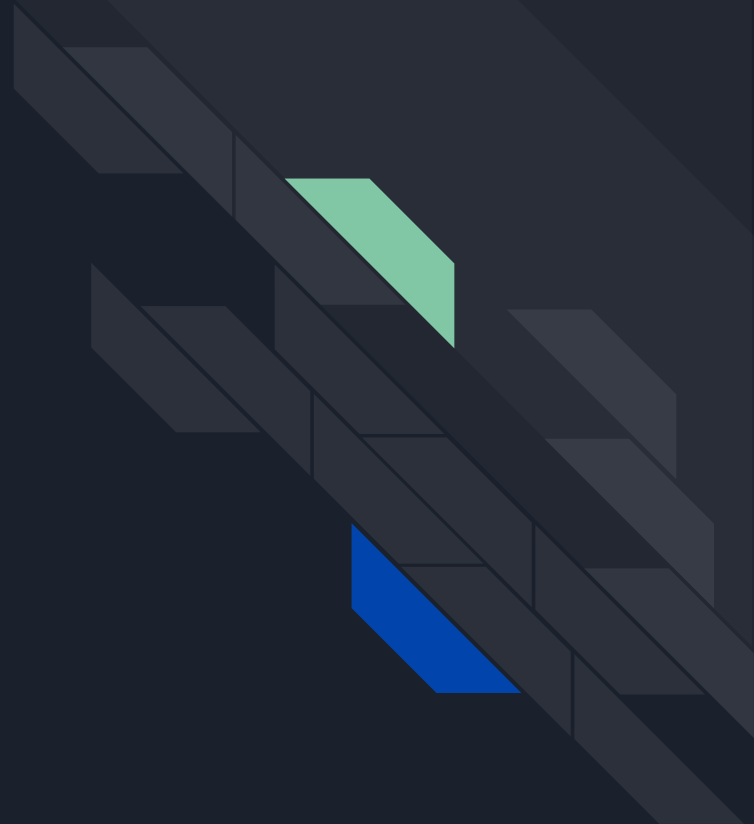
# Talk Outline

- Linear Algebra
    - Vectors
    - Vector Spaces
    - Linearity
- Machine Learning
    - Training
    - Single Layers
    - Neural Networks
- NLP & Word Embeddings
    - Character Encodings
    - Higher Dimensionality
    - Encoders and Decoders

# Talk Outline (Cont'd)

- Evaluating Word Embeddings using Probing
    - A Different Result
- Criticisms of Probing
    - No Control
    - Model Variety
    - Correlation and Causation
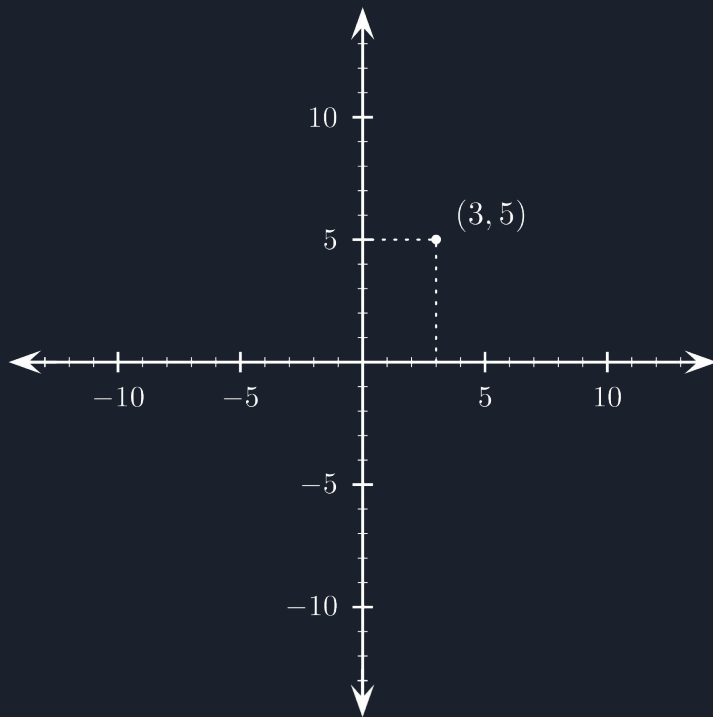- Experiments on Large Language Models

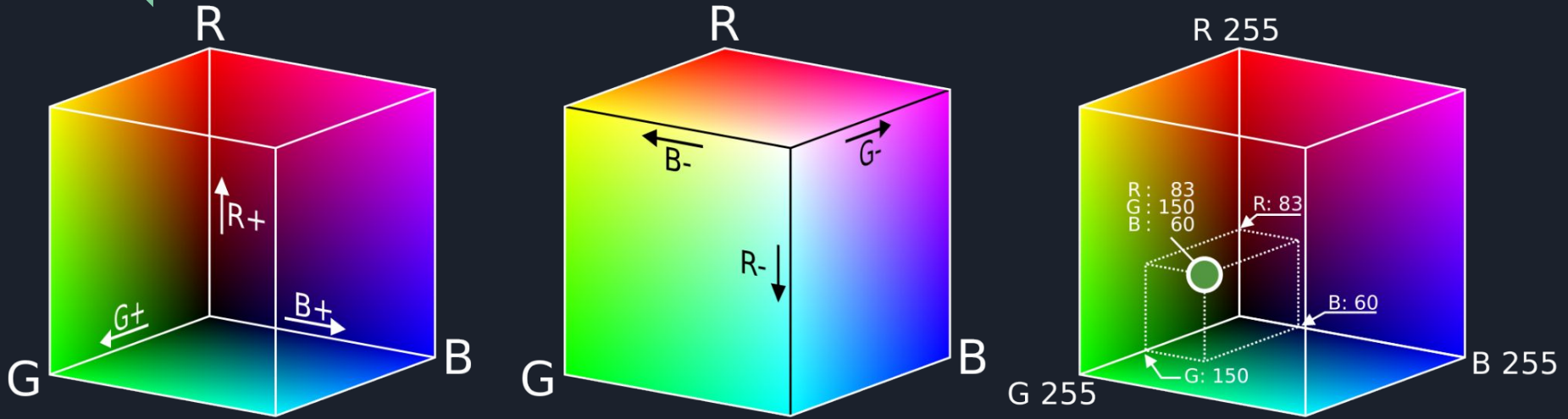# Linear Algebra

# Background - Vectors and Vector Spaces

- We can think of vectors as lists of numbers

$$\begin{bmatrix} 0.24 \\ 8.02 \\ -3.4 \\ 3.14 \end{bmatrix}$$

# Background - Vectors and Vector Spaces

# Background - Vectors and Vector Spaces
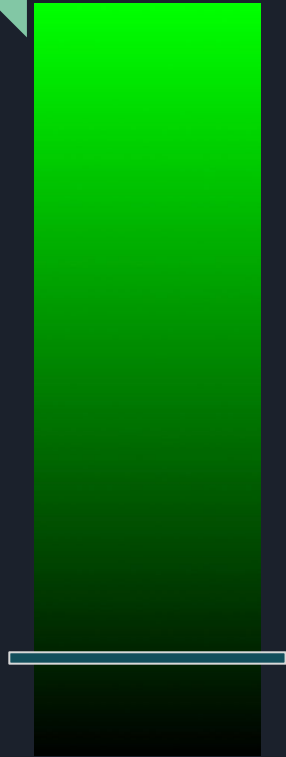


Adapted From: WikiMedia

# Background - Linear Transformations

- Treating dimensions individually
- Combining dimensions individually by summing
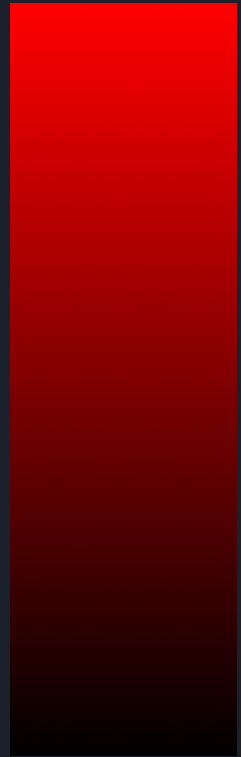
# Background - Linear Transformations

# Background - Linear Transformations

# Background - Linear Transformations

# Background - Linear Transformations

# Background - Linear Transformations

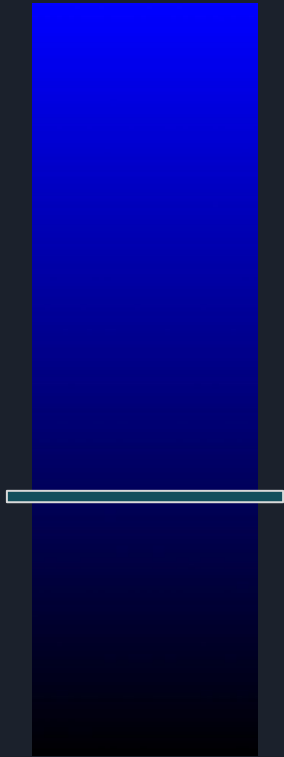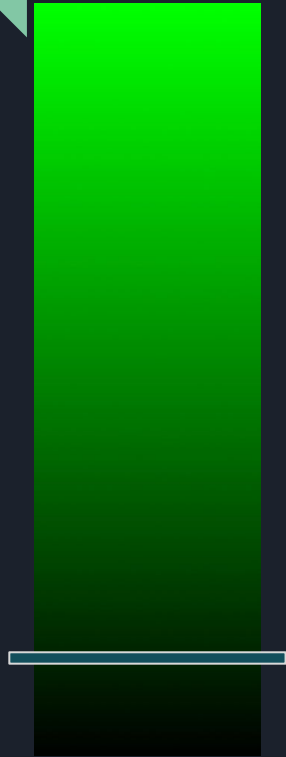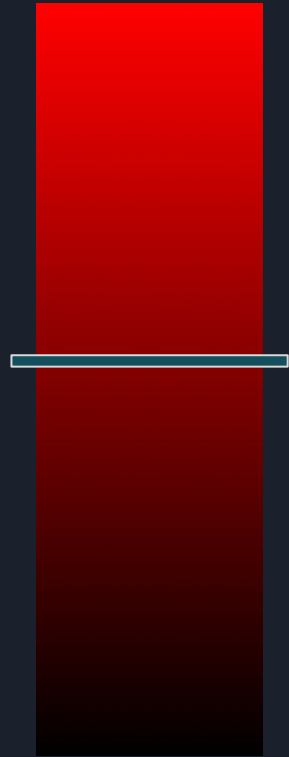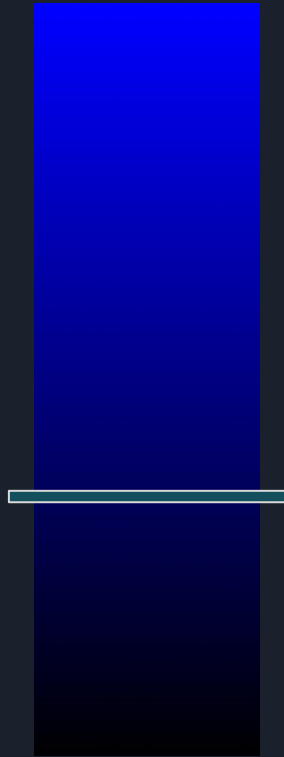# Background - Linear Transformations

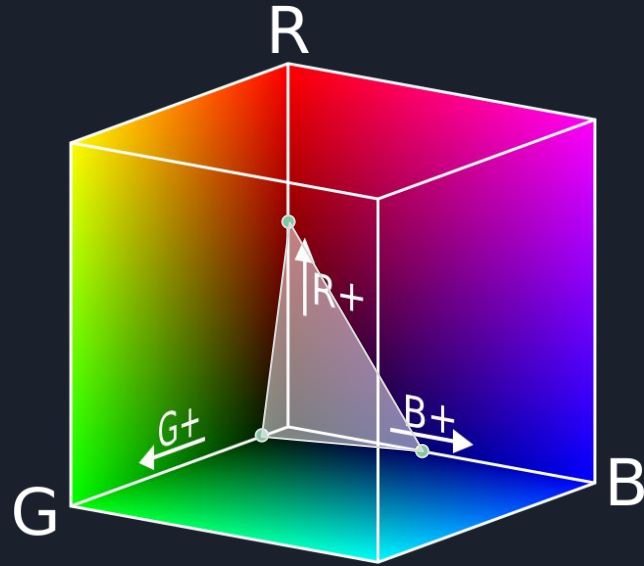Luma(Blue) = 1

Luma(Red) = 3

Luma(Green) = 4

# Background - Linear Transformations

Luma(Blue) = 1                    Luma(Blue) = 1/8

Luma(Red) = 3                     Luma(Red) = 3/8

Luma(Green) = 4                   Luma(Green) = 4/8

# Background - Linear Transformations

Luma(r·Red + g·Green + b·Blue)

= ⅜·r + ½·g + ⅛·b

# Background - Matrix Multiplication

$$\begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix} \begin{bmatrix} A \\ B \\ C \end{bmatrix}$$

$$\begin{bmatrix} A \\ B \\ C \\ D \\ E \end{bmatrix} \longrightarrow \begin{bmatrix} a \\ b \end{bmatrix}$$

# Machine Learning

$$Parameters$$

$$\downarrow$$

$$Input \longrightarrow [Model] \longrightarrow Output$$

# Background - Machine Learning

- Trained by providing pairs of input and output

# Background - Machine Learning

- Trained by providing pairs of input and output
  - Apply model to the input

# Background - Machine Learning

- Trained by providing pairs of input and output
    - Apply model to the input
    - Compare the output with the expected output

# Background - Machine Learning

- Trained by providing pairs of input and output
    - Apply model to the input
    - Compare the output with the expected output
    - Use that information to update parameters

$$\text{Parameters} \longleftarrow \text{Feedback}$$

$$\downarrow \qquad\qquad\qquad \uparrow$$

$$\text{Input} \longrightarrow [\text{Model}] \longrightarrow \text{Output}$$

# Background - Single Layer

- Matrix Multiplication
- Linear Classifier

# Background - Single Layer

## Linear Separation

# Background - Neural Networks

## Activation Functions

# Background - Neural Networks

$$\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} \begin{bmatrix} \ \end{bmatrix} \sigma \atop \longrightarrow \begin{bmatrix} a \\ b \\ c \\ d \\ e \\ f \end{bmatrix} \begin{bmatrix} \ \end{bmatrix} \sigma \atop \longrightarrow \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} \begin{bmatrix} \ \end{bmatrix} \sigma \atop \longrightarrow \begin{bmatrix} a \\ b \end{bmatrix}$$

Natural Language Processing

# NLP - Word Embeddings

How do we represent a word as a vector?

- Character Encodings?
- Many, Many Dimensions
- Encoders and Decoders

# NLP - Character Encoding?



**Figure 1.** ASCII-based embedding of the word "vector"

# NLP - Character Encoding?

Limitations:

- Limited in length
- Doesn't play well as a vector

# NLP - A Dimension for Every Word



**Figure 2.** Simple embedding for the word "vector"

# NLP - A Dimension for Every Word

Advantages:

-   Plenty of space for machine learning

Disadvantages:

-   1.3 Million dimensions

# NLP - A Useless Transformation

$$\text{Parameters}$$

$$\downarrow$$

$$\text{Input} \longrightarrow \text{[Model]} \longrightarrow \text{Output}$$

# NLP - A Useless Transformation

Parameters

"Vector"

Input $\longrightarrow$ [Model] $\longrightarrow$ Output
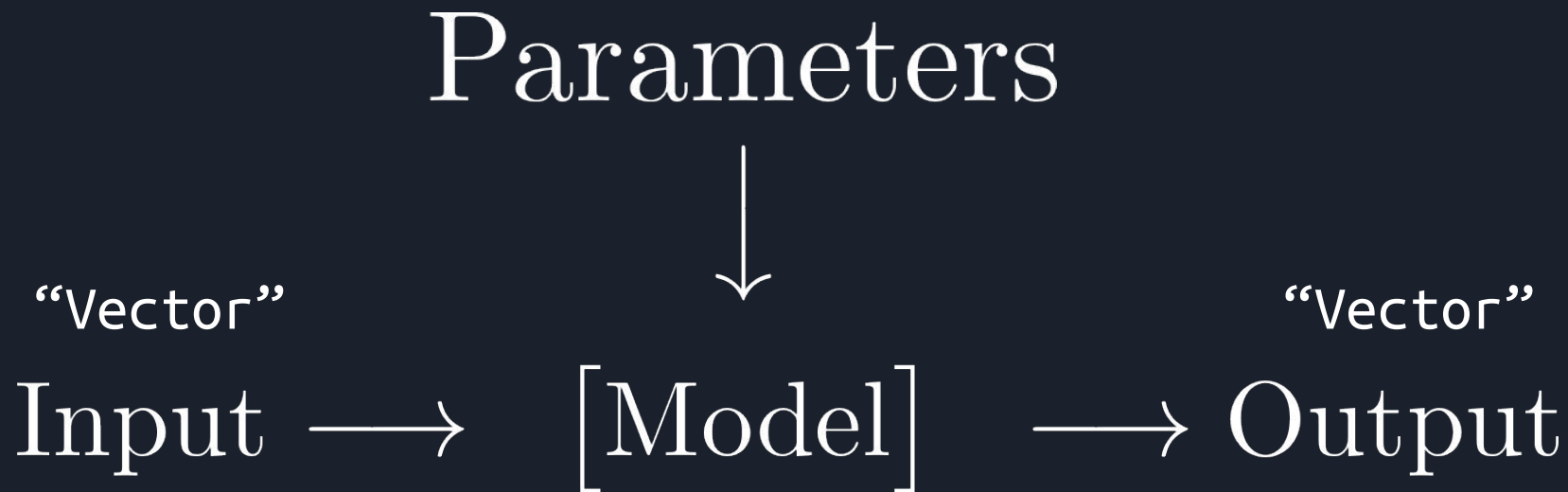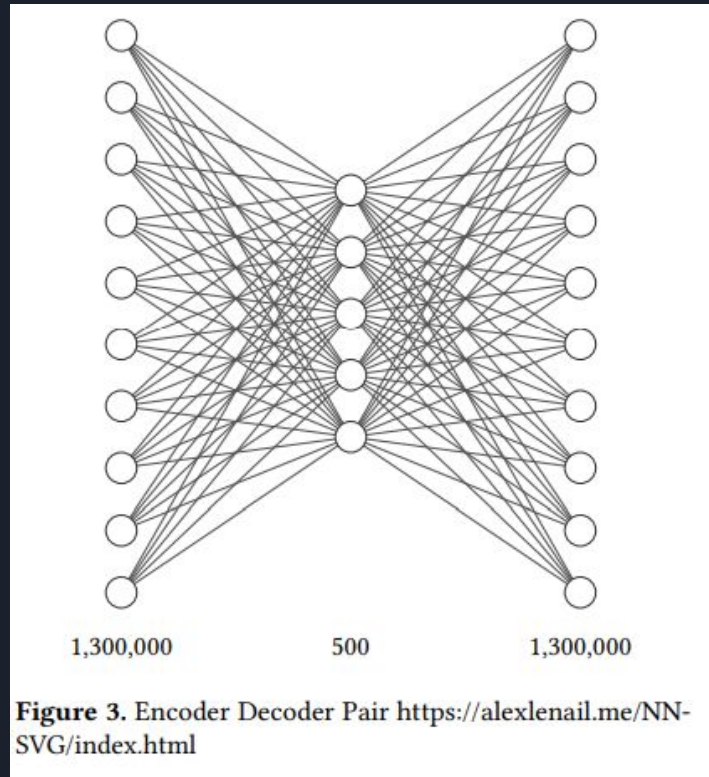
# NLP - A Useless Transformation

Parameters

$\downarrow$

"Vector"

"Vector"

Input $\longrightarrow$ [Model] $\longrightarrow$ Output

# NLP - Encoders and Decoders



1,300,000      500      1,300,000

**Figure 3.** Encoder Decoder Pair https://alexlenail.me/NN-SVG/index.html

Parameters

↓

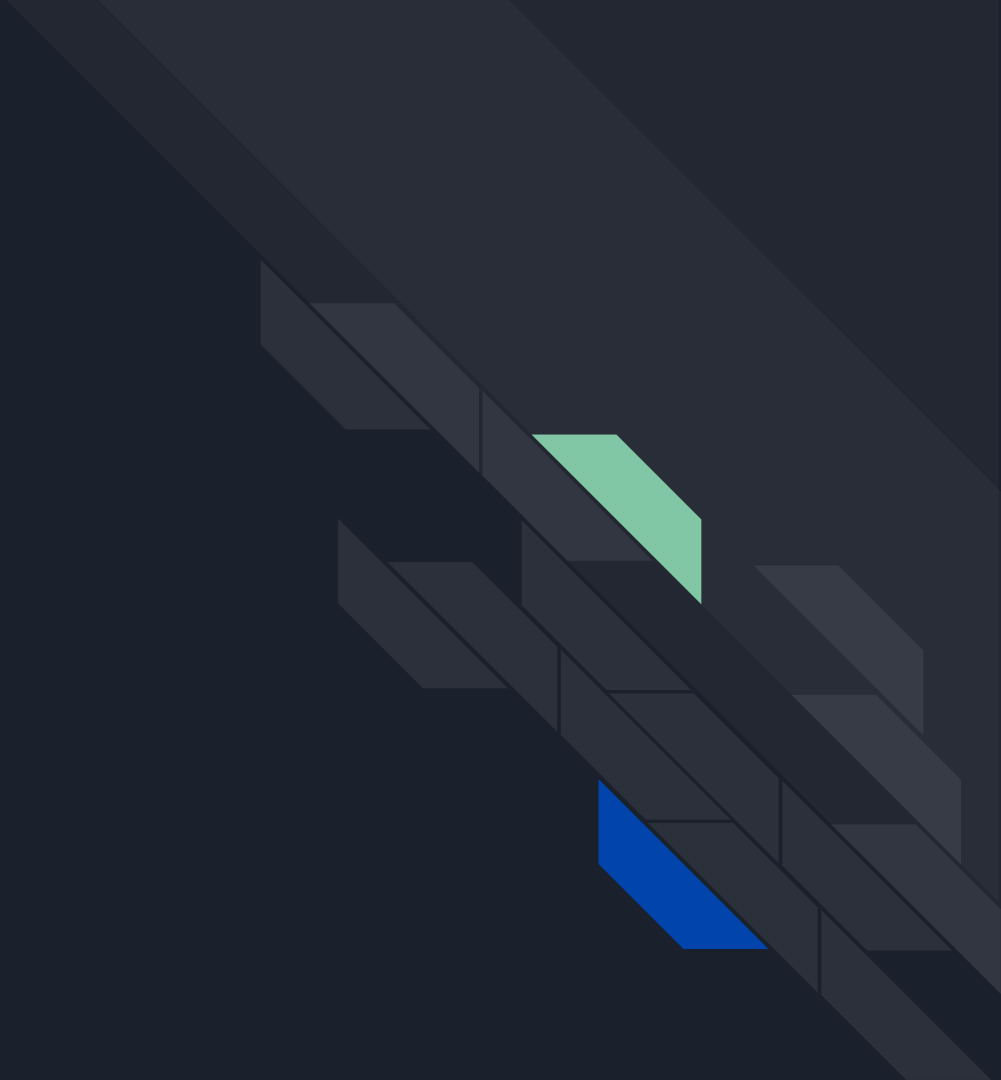Input ⟶ [Model] ⟶ Output

# Probing

# Probing - Assessing Word Embeddings

- Part of Speech?

- Plurality?

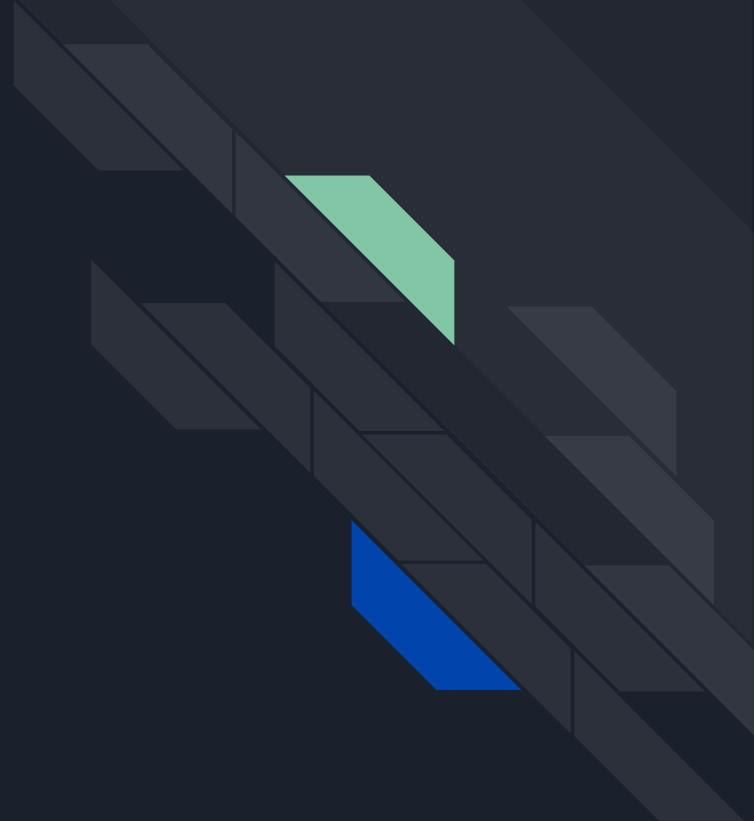- Ends-in-'s'-ness?

# Probing - A Different Application

"Traditional" Machine Learning:

The parameters are the result, the accuracy is just a measurement

Probing:

The accuracy in training is the result, the parameters are coincidental

# Criticisms

# Criticisms

No Control Reference

Models Can Vary

Correlation and Causation

# Criticisms - No Control

What is the baseline?

# Criticisms - Model Variety

When does one use a Linear Classifier?

What activation function should be used if any?

# Criticisms - Correlation and Causation
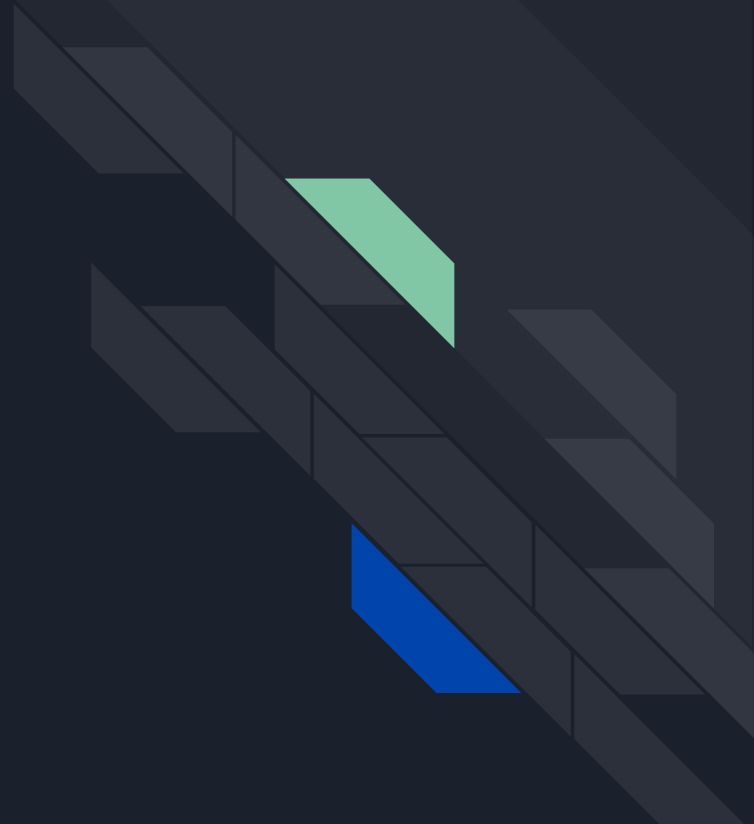
How does one remove a property?

# Criticisms - Correlation and Causation

How do we know the property we want is present?

- Retrain the encoder
- Modifying the embedded vector

# Using Our Toolkit

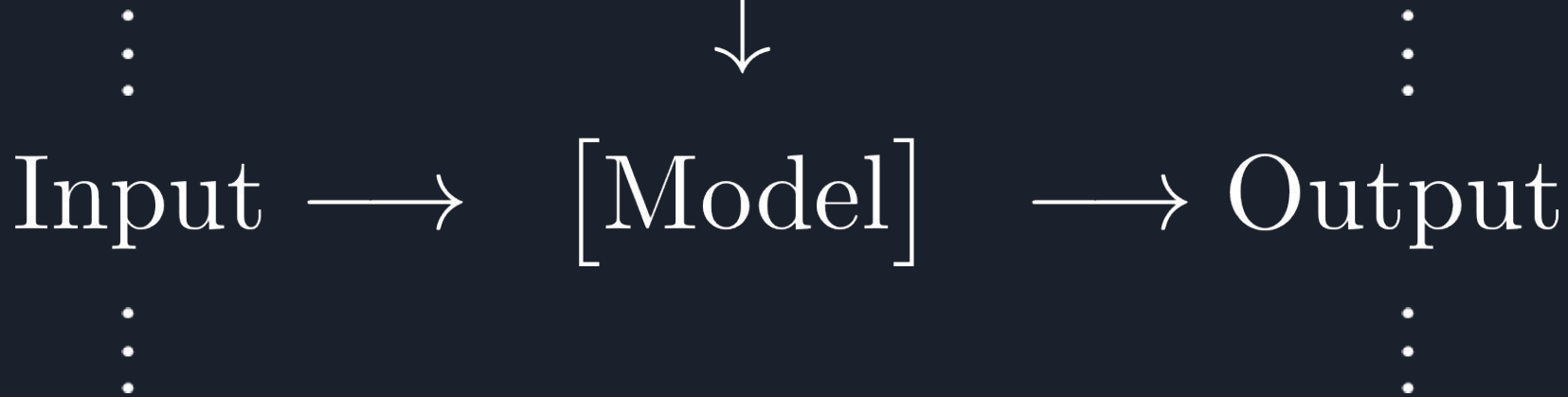# Experiments on Large Language Models

$$Parameters$$

$$\downarrow$$

$$Input \longrightarrow [Model] \longrightarrow Output$$

# Experiments on Large Language Models

Parameters

$\downarrow$

Input $\longrightarrow$ [Model] $\longrightarrow$ Output

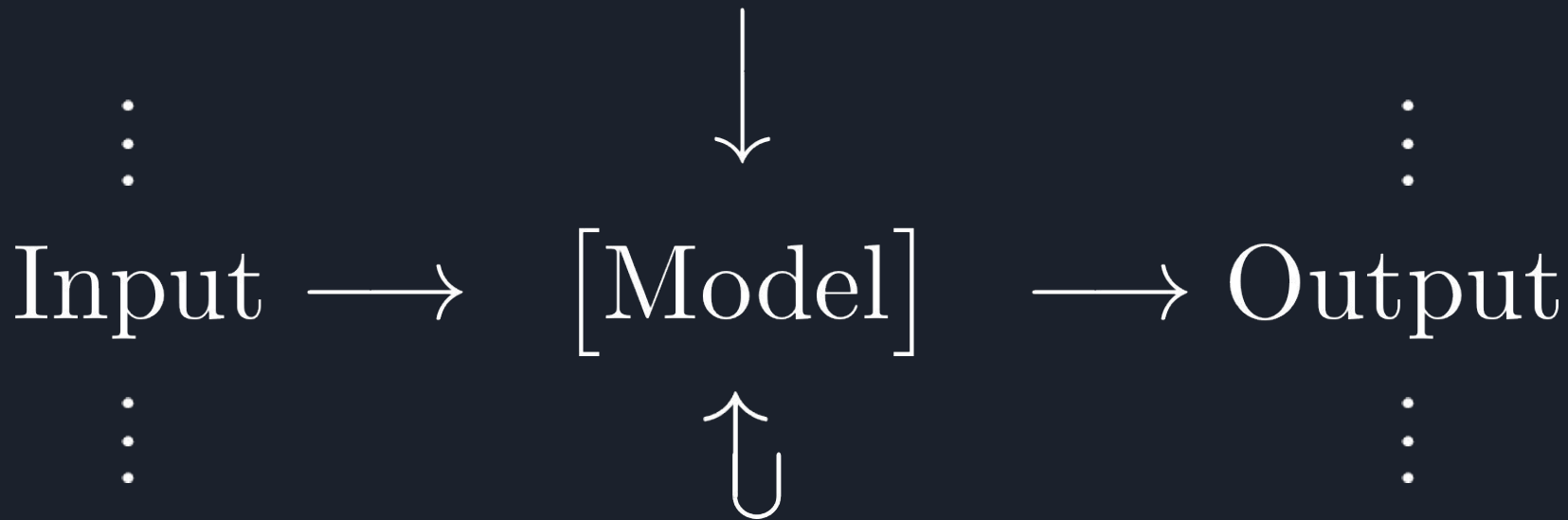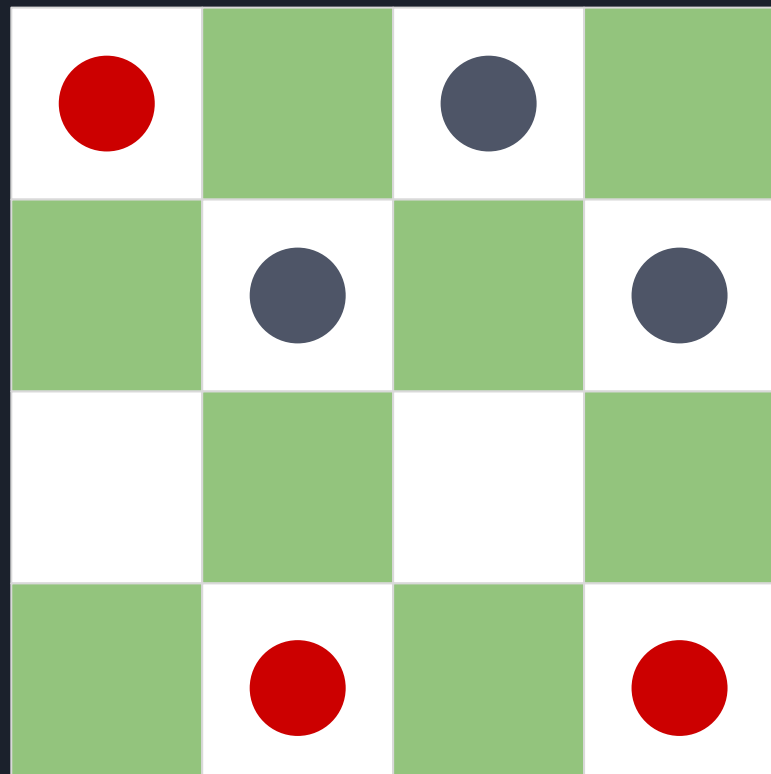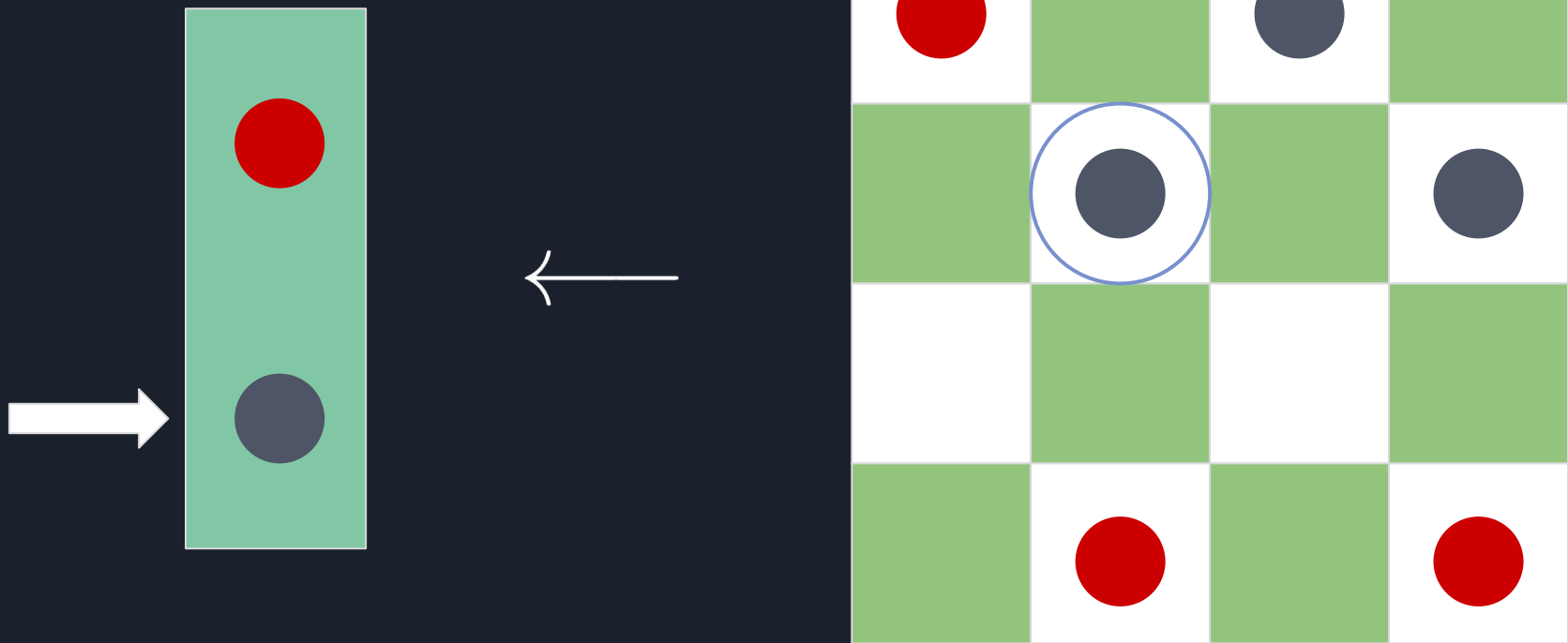# Experiments on Large Language Models

# Experiments on Large Language Models

# Experiments on Large Language Models

# Experiments on Large Language Models

# Experiments on Large Language Models
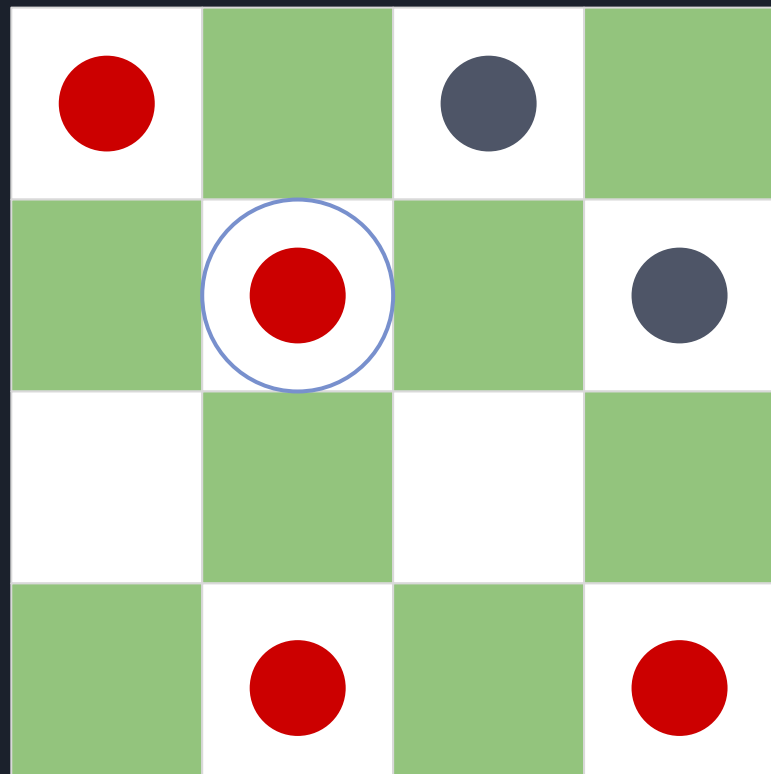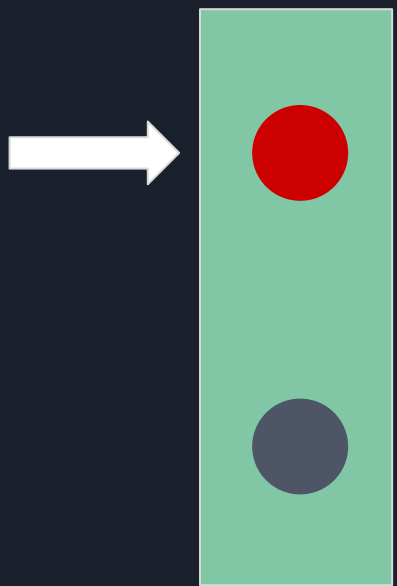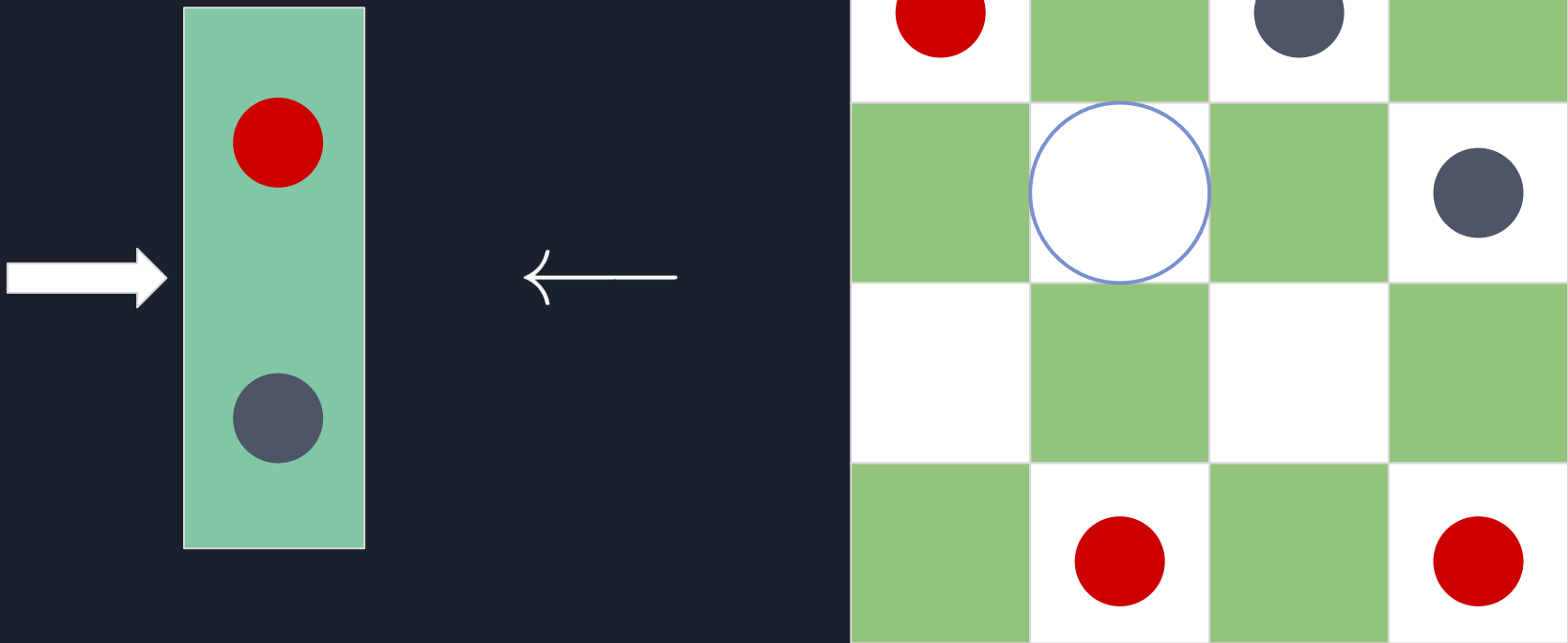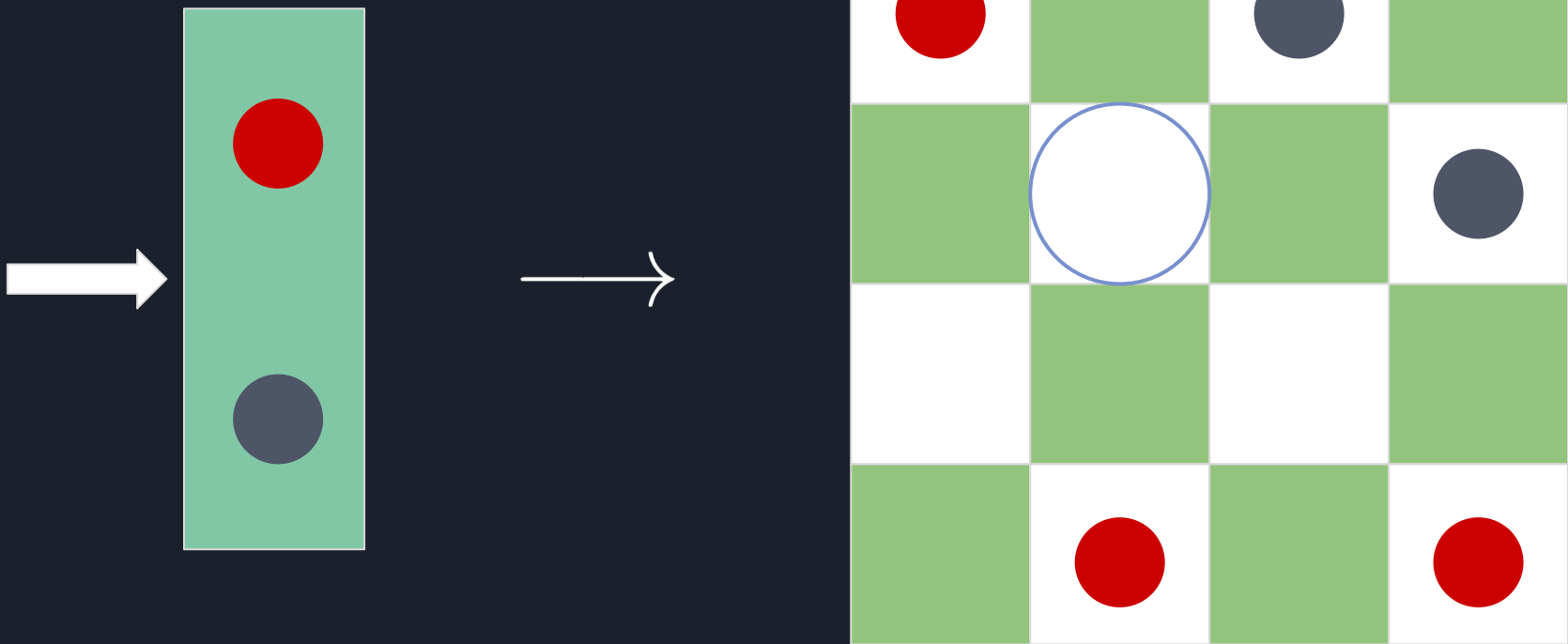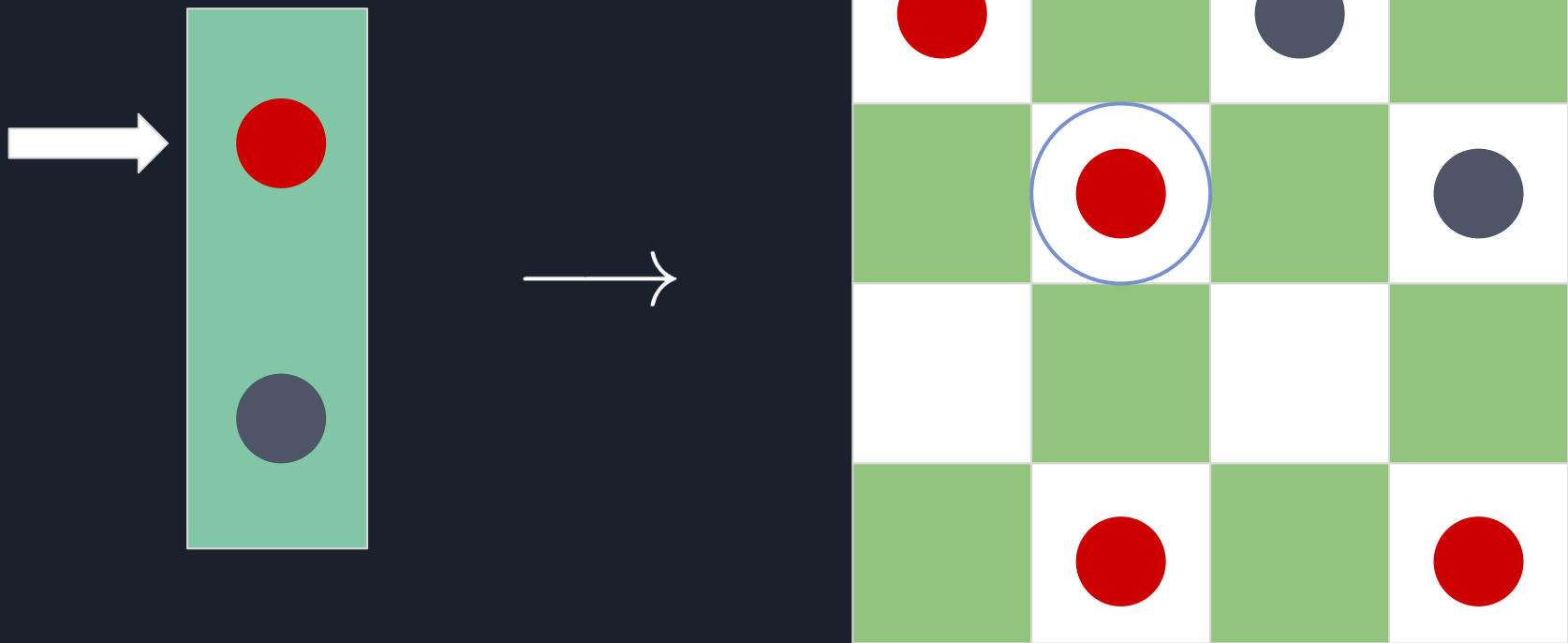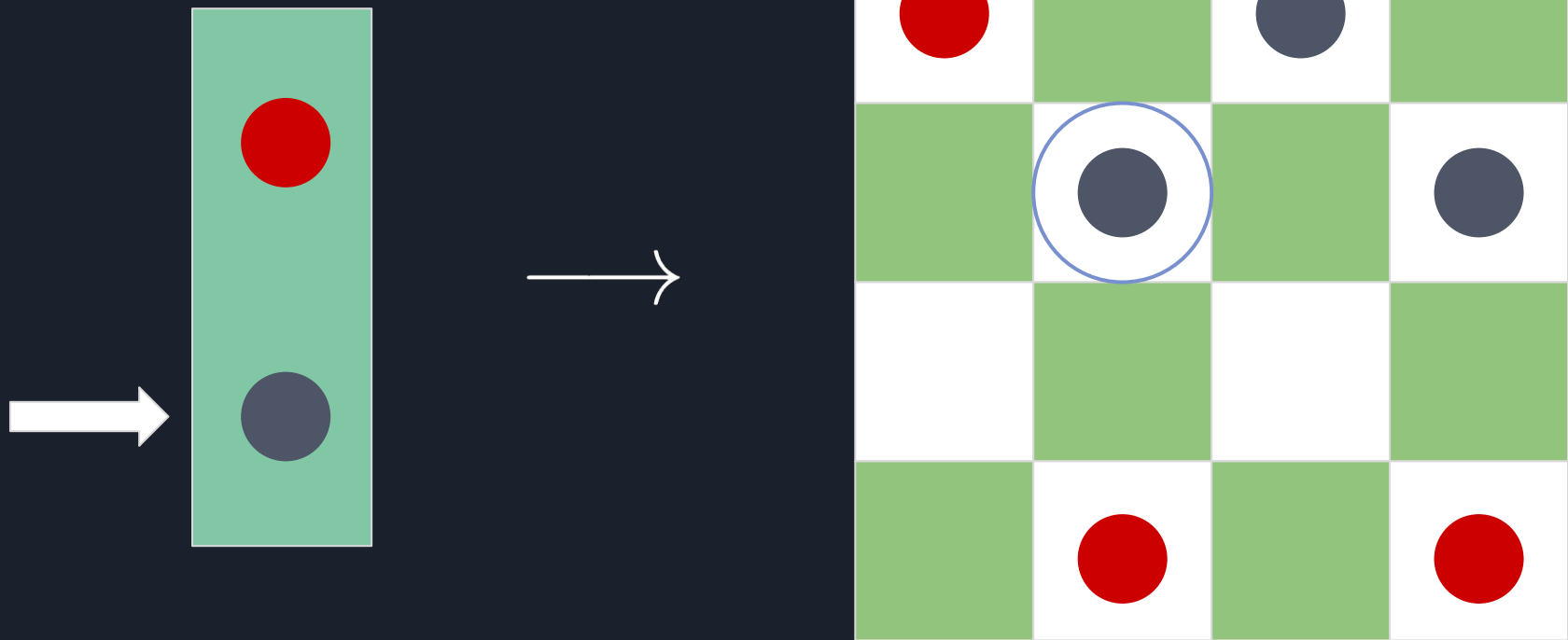
# Experiments on Large Language Models

# Experiments on Large Language Models

# Experiments on Large Language Models

# Conclusion

# Conclusion

- Linear Algebra
- Machine Learning
- NLP & Word Embeddings
- Evaluating Word Embeddings using Probing
- Criticisms of Probing
- Experiments on Large Language Models

# References

[1]  Guillaume Alain and Yoshua Bengio. 2016. Understanding intermediate layers using linear classifier probes. (2016). https://doi.org/10.48550/ ARXIV.1610.01644

[2]  Yonatan Belinkov. 2022. Probing Classifiers: Promises, Shortcomings, and Advances. Computational Linguistics 48, 1 (04 2022), 207–219. https://doi.org/10.1162/coli_a_00422 arXiv: https://direct.mit.edu/coli/article-pdf/48/1/207/2006605/coli_a_00422.pdf

[3]  Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single $&!#* vector: Probing sentence embeddings for linguistic properties. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Melbourne, Australia, 2126–2136. https://doi.org/10.18653/v1/P18-1198

# References

[4]  Arne Köhn. 2015. What's in an Embedding? Analyzing Word Embeddings through Multilingual Evaluation. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, Lisbon, Portugal, 2067–2073. https://doi.org/10.18653/v1/D15-1246

[5]  Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2022. Emergent World Representations: Exploring a Sequence Model Trained on a Synthetic Task. (2022). https://doi.org/10.48550/ARXIV.2210.13382

[6]  Wiktionary. 2023. Statistics — Wiktionary. https://en.wiktionary.org/wiki/Special:Statistics. [Online; accessed 01-March-2023].

# Questions?