# Lesson 5: Generating and Encoding Text with Transformers

Pearson

# 5.1

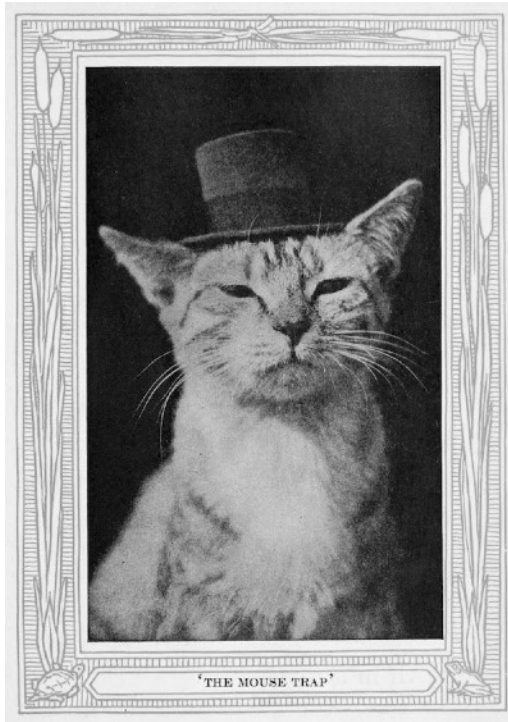## The Natural Language Processing Pipeline

# Natural Language Processing (NLP)

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet. Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

$\longrightarrow$

[1, 3, 1, 1, 2, 0, 1, 0]
[0, 1, 4, 0, 0, 1, 1, 1]
[3, 0, 1, 1, 2, 2, 3, 2]
[0, 1, 1, 1, 0, 3, 2, 3]
[1, 2, 1, 2, 2, 0, 0, 0]
[1, 0, 1, 1, 0, 1, 1, 1]
[0, 2, 0, 0, 2, 2, 0, 0]
[1, 1, 1, 1, 0, 1, 1, 1]

Pearson

# Image Vectorization



| 245 | 238 | 222 | 255 |
|-----|-----|-----|-----|
| 233 | 0 | 17 | 254 |
| 255 | 6 | 3 | 223 |
| 250 | 9 | 11 | 242 |
| 251 | 247 | 245 | 232 |

# NLP Pipeline

# NLP Pipeline

| Tokenization | → | Preprocessing | → | Feature Extraction |
|---|---|---|---|---|

- Sentence level
- Word level
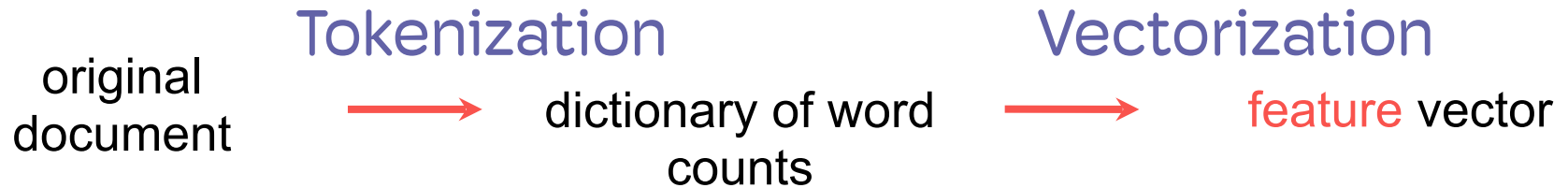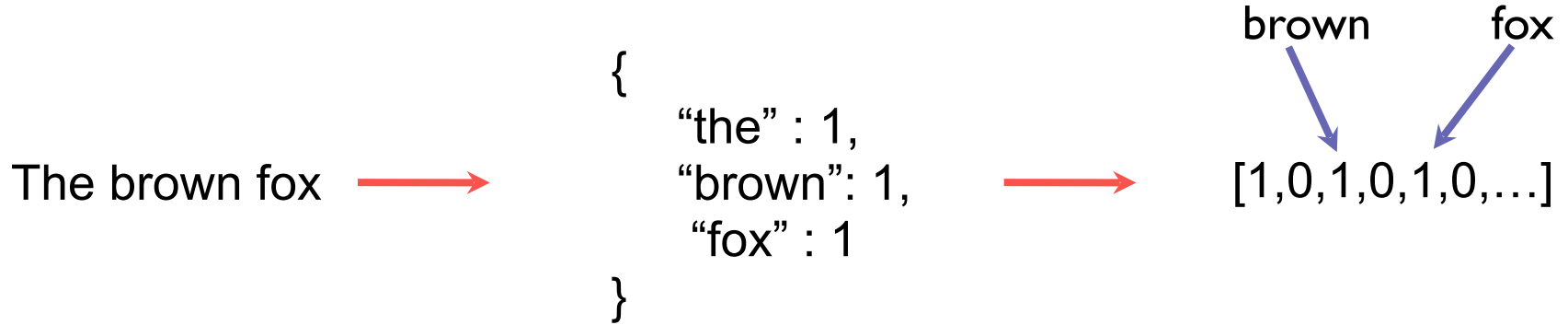- Character level
- Byte pair encoding

- Punctuation filtering
- Stop word removal
- Stemming
- Lemmatization

- Bag of words
- TF-IDF
- POS tagging
- word2vec

# Terminology

- **Document:** Single row of data/corpus

- **Corpus:** Entire set of all documents

- **Vocabulary:** Set of all words in corpus

- **Vector:** Mathematical representation of document (counts of word occurrences)

# Bag of Words

The brown fox $\longrightarrow$

```
{
    "the" : 1,
    "brown": 1,
    "fox" : 1
}
```

$\longrightarrow$

brown    fox

[1,0,1,0,1,0,…]

original document

**Tokenization**

$\longrightarrow$ dictionary of word counts

**Vectorization**

$\longrightarrow$ feature vector

# Bag of Words (Bernoulli)

| | red | brown | jumps | the | fox | panda |
|---|---|---|---|---|---|---|
| doc0 | 0 | 1 | 0 | 1 | 1 | 0 |
| doc1 | 1 | 0 | 0 | 1 | 1 | 0 |
| doc2 | 1 | 0 | 0 | 1 | 0 | 1 |
| doc3 | 0 | 0 | 1 | 1 | 1 | 0 |

# Bag of Words (Multinomial)

| | red | brown | jumps | the | fox | panda |
|---|---|---|---|---|---|---|
| doc0 | 0 | 2 | 0 | 4 | 2 | 0 |
| doc1 | 1 | 0 | 0 | 1 | 1 | 0 |
| doc2 | 2 | 0 | 0 | 2 | 0 | 2 |
| doc3 | 0 | 0 | 1 | 2 | 2 | 0 |

# Dedicated NLP Libraries

- spaCy: https://spacy.io

- fastText: https://fasttext.cc

- Gensim: https://radimrehurek.com/gensim/

- AllenNLP: https://allenai.org/allennlp

- flair: https://github.com/flairNLP/flair

- fairseq: https://github.com/facebookresearch/fairseq

# 5.2

## Generative Models of Language

# Probabilistic Model of Natural Language

$$P(w_1, \ldots, w_m) = \prod_{t=1}^{m} P(w_t \mid w_{1:t-1})$$

Pearson

# Language Models

History of words
that came before

$$P(w_t \mid w_{1:t-1})$$

Probability of
next word

# Language Models

Often not "full" history (infinite context)

$$P(w_t \mid w_{t-5:t-1})$$

Probability of next word

Pearson

# Language Models

- **n-gram model:** Fixed window of previous $n$ words

- **Neural/RNN:** Learned embeddings of $n$ words

- **Large language model (LLM):** Large scale self- and semi-supervised pretraining of bidirectional models

# Causal Language Modeling

# 5.3

**Generating Text with Transformers Pipelines**

# Live Coding

# 5.4

## Deconstructing Transformers Pipelines
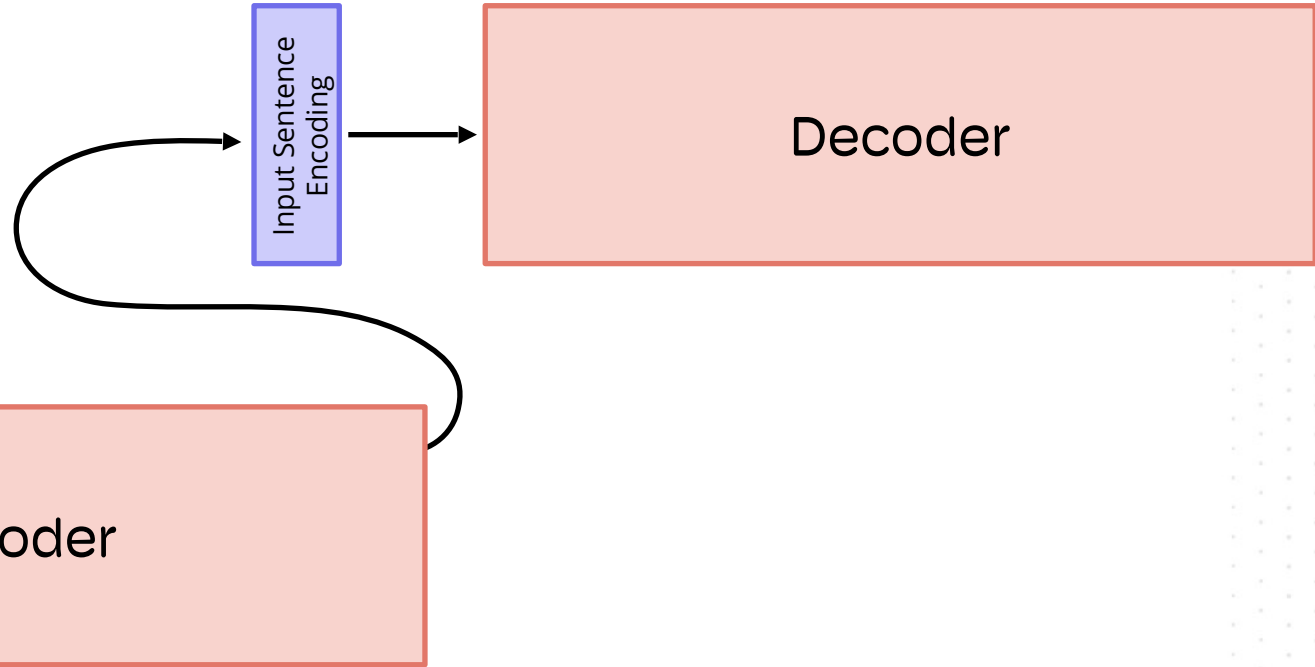
# Live Coding

# 5.5

## Decoding Strategies

# Live Coding

# 5.6

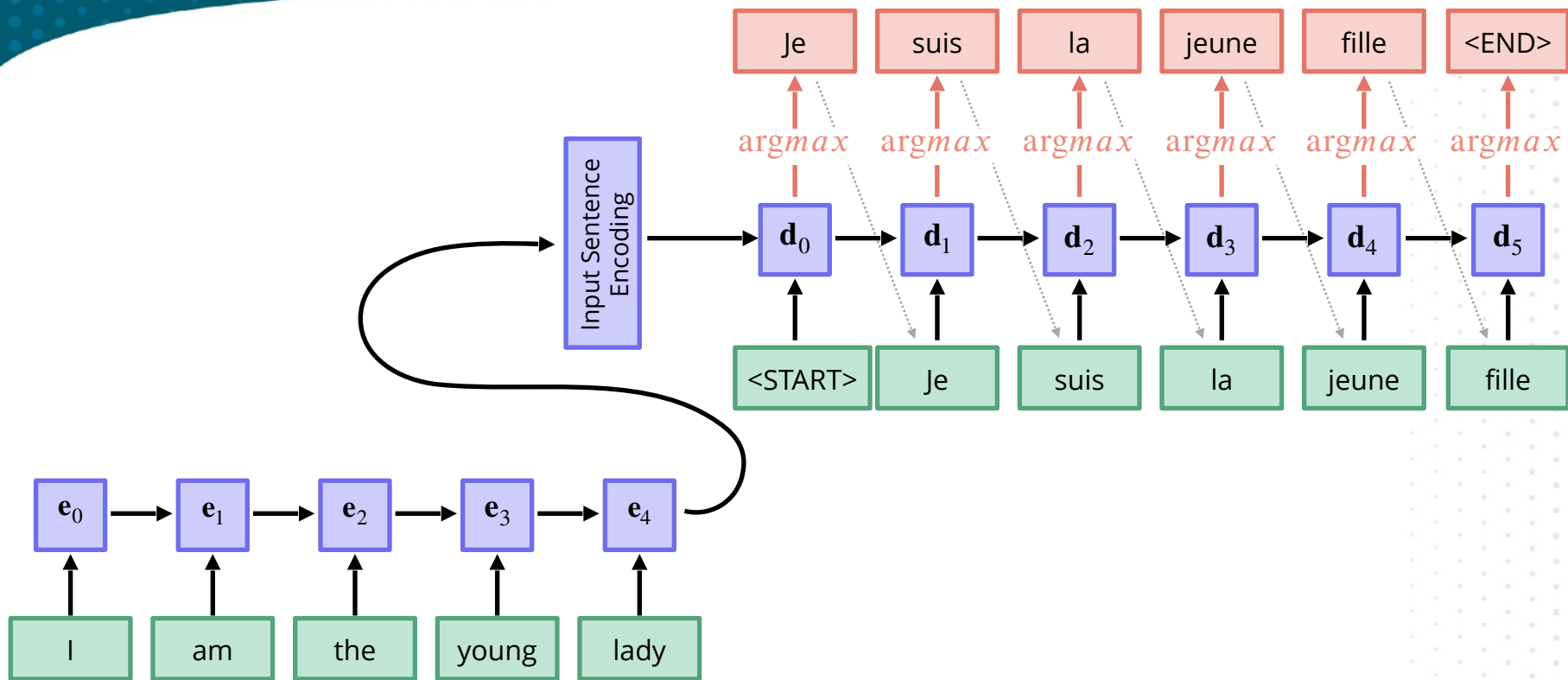## Transformers are Just Latent Variable Models for Sequences
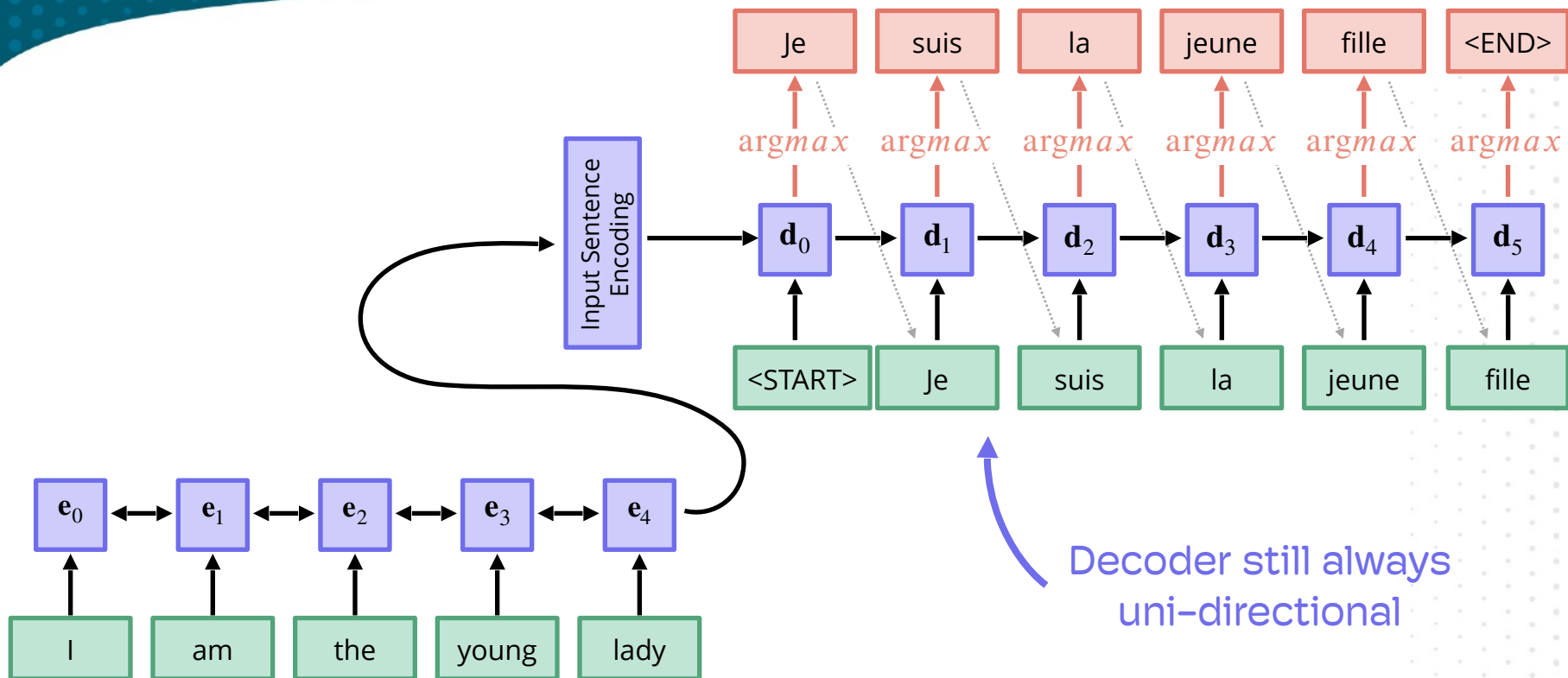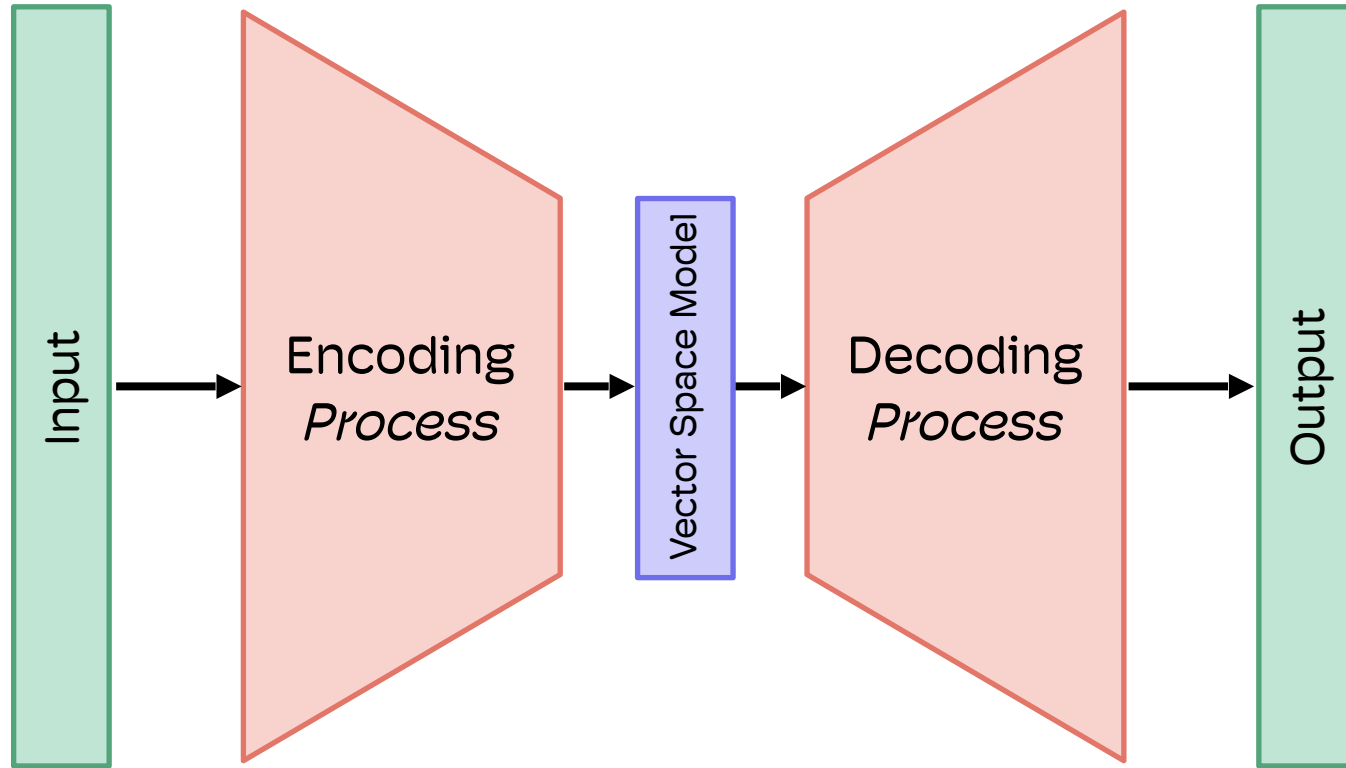
# Text-to-Text (Translation, Summarization, Other)

Text

Images

Audio

Video

Probably a Transformer

Text

Images

Audio

Video

Pearson

# seq2seq

# seq2seq

# Bidirectional seq2seq



Decoder still always uni-directional

# Encoding Natural Language



Input → Encoding *Process* → Vector Space Model → Decoding *Process* → Output

# Transformer

I

am

the

young

lady

FF Neural Network

FF Neural Network

FF Neural Network

Internal States

**Can "see" the entire input**

Pearson

# Encoder Only Transformer



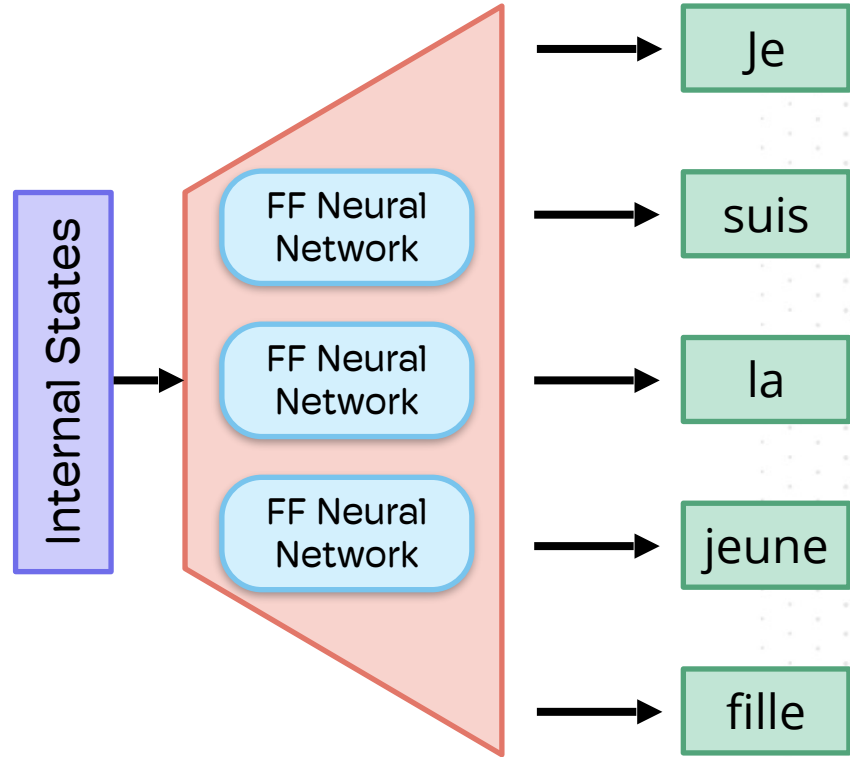| I |
| am |
| the |
| young |
| lady |

FF Neural Network

FF Neural Network

FF Neural Network

Internal States

**Typically for downstream tasks (i.e. BERT)**

Pearson

# Decoder Only Transformer

**Only has access to the previous words**



Internal States

FF Neural Network

FF Neural Network

FF Neural Network

Je

suis

la

jeune

fille

Pearson

# Decoder Only Transformer

**Typically for generation (i.e. GPT)**



Internal States

FF Neural Network → Je

FF Neural Network → suis

FF Neural Network → la

→ jeune

→ fille

Pearson

# 5.7

## Visualizing and Understanding Attention

# Live Coding

# 5.8

## Turning Words into Vectors

# Encoding Natural Language



Input → Encoding *Process* → Vector Space Model → Decoding *Process* → Output

# Natural Image Manifold

# Natural Language Manifold



Mona Lisa

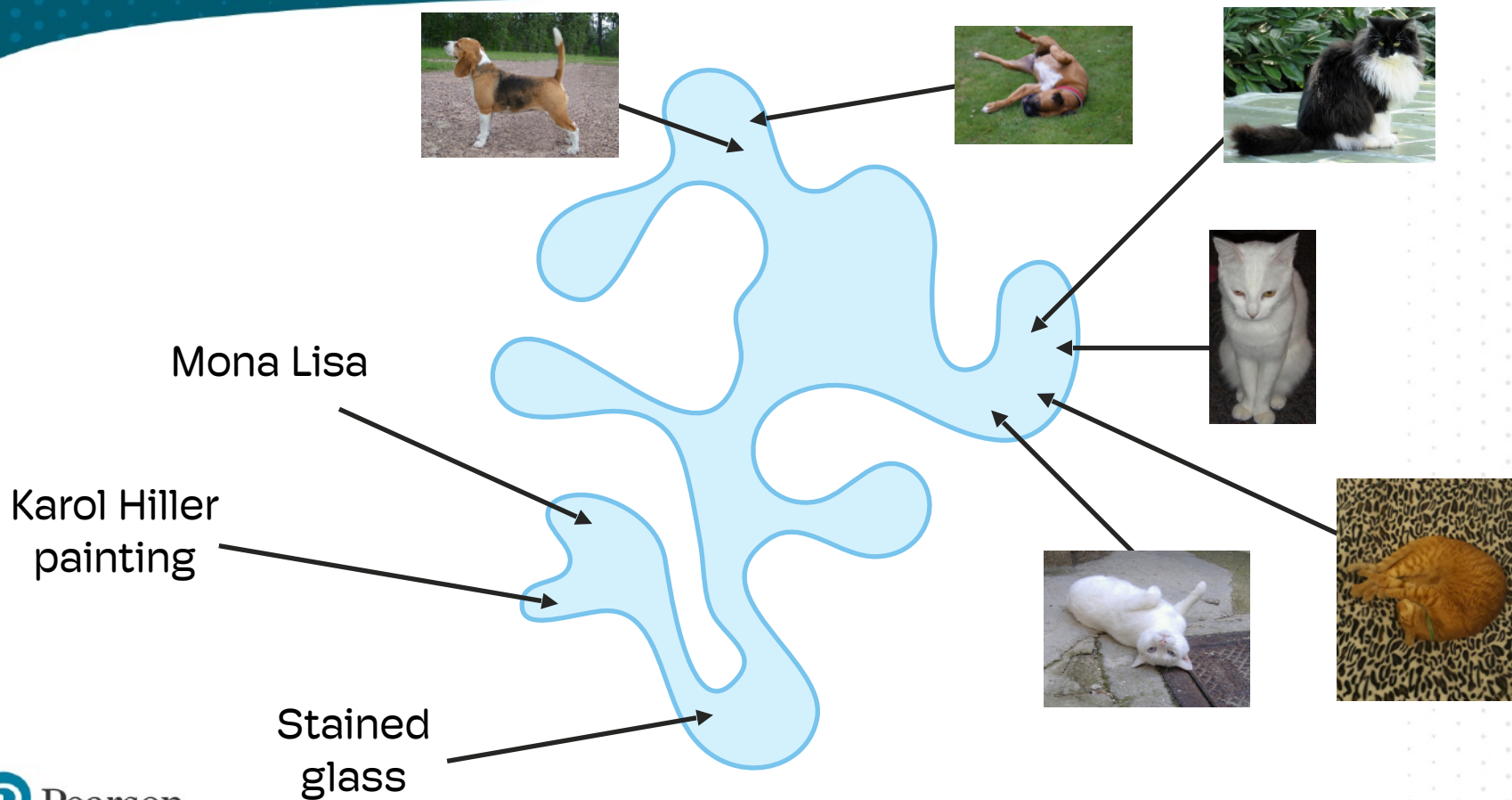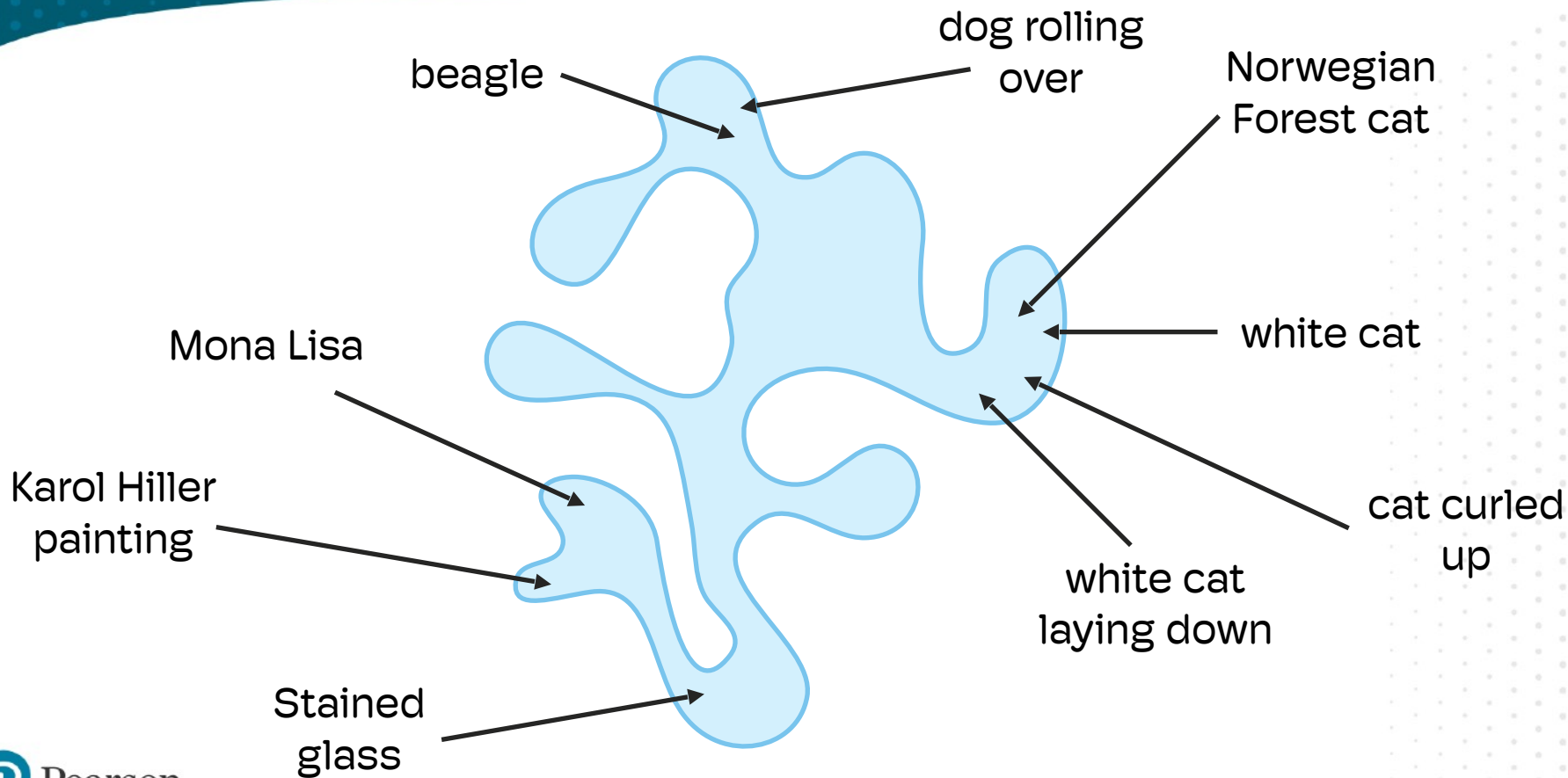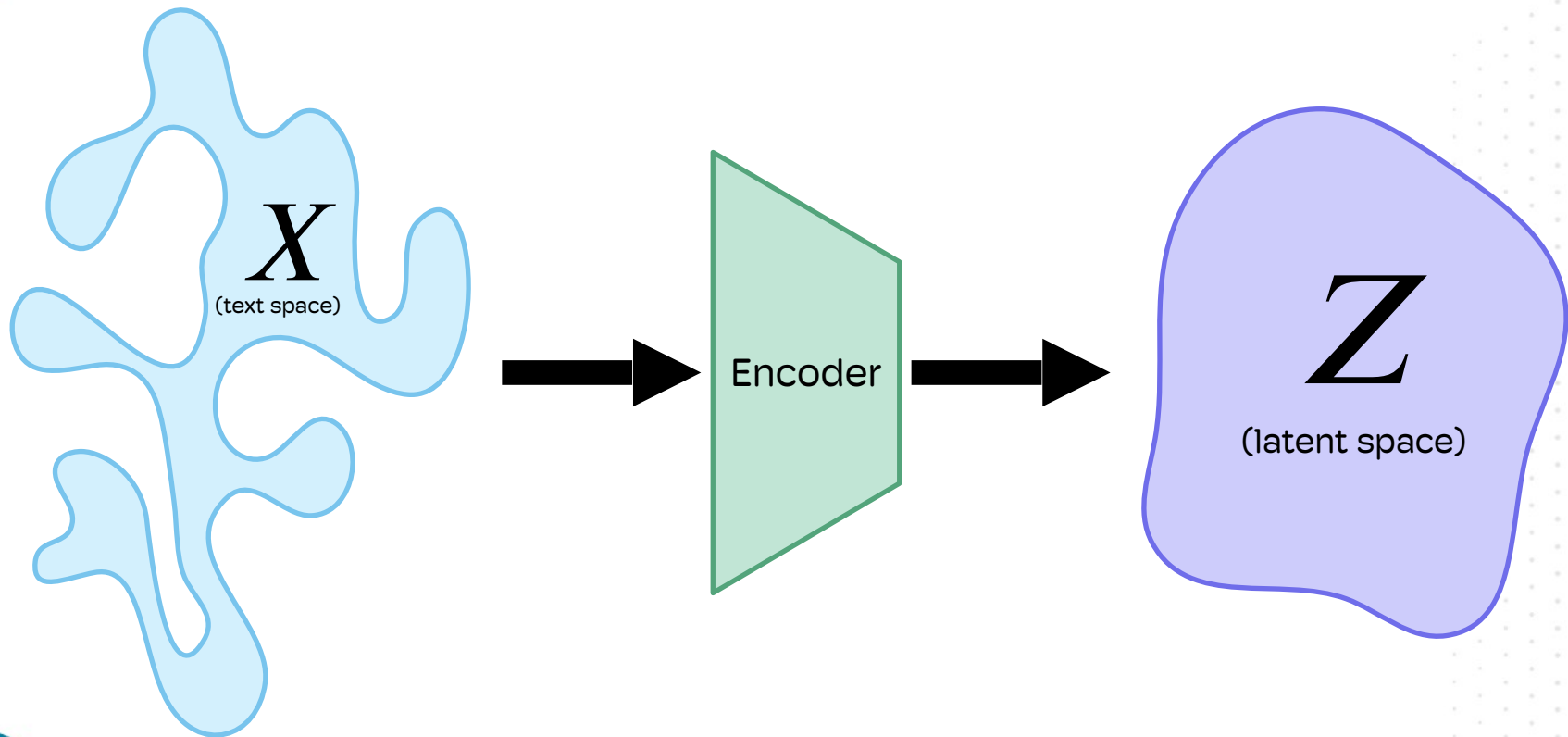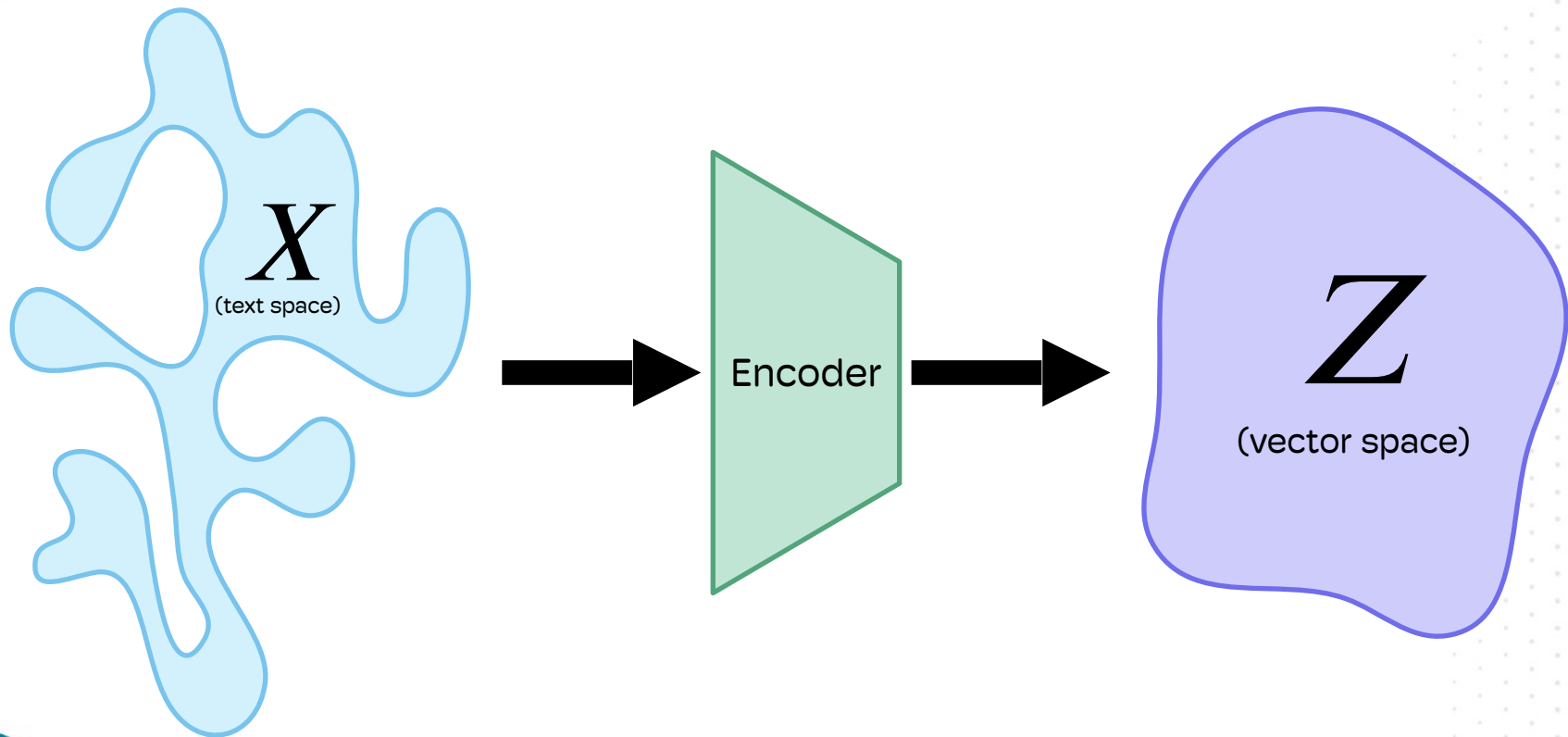Karol Hiller painting

Stained glass

# Natural Language Manifold
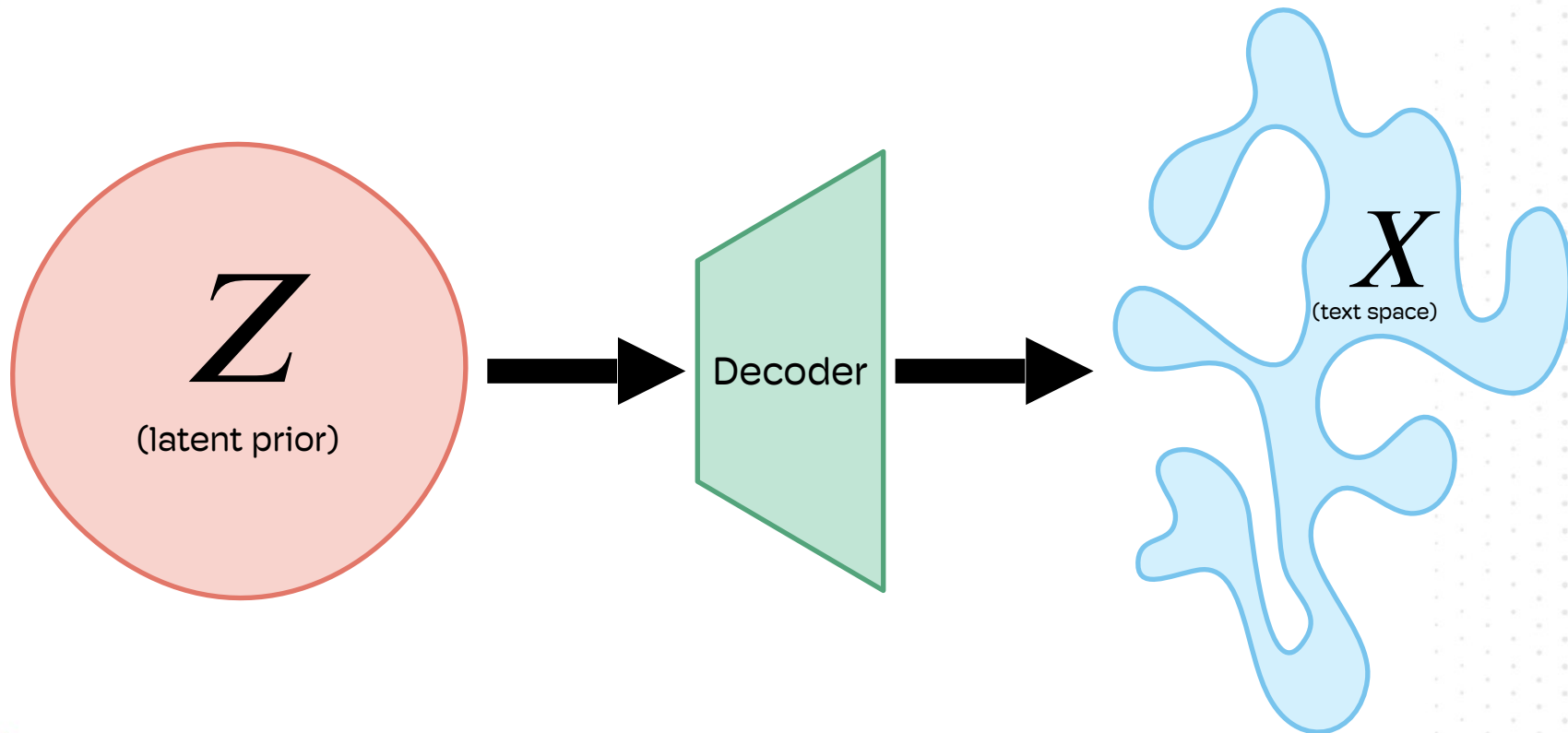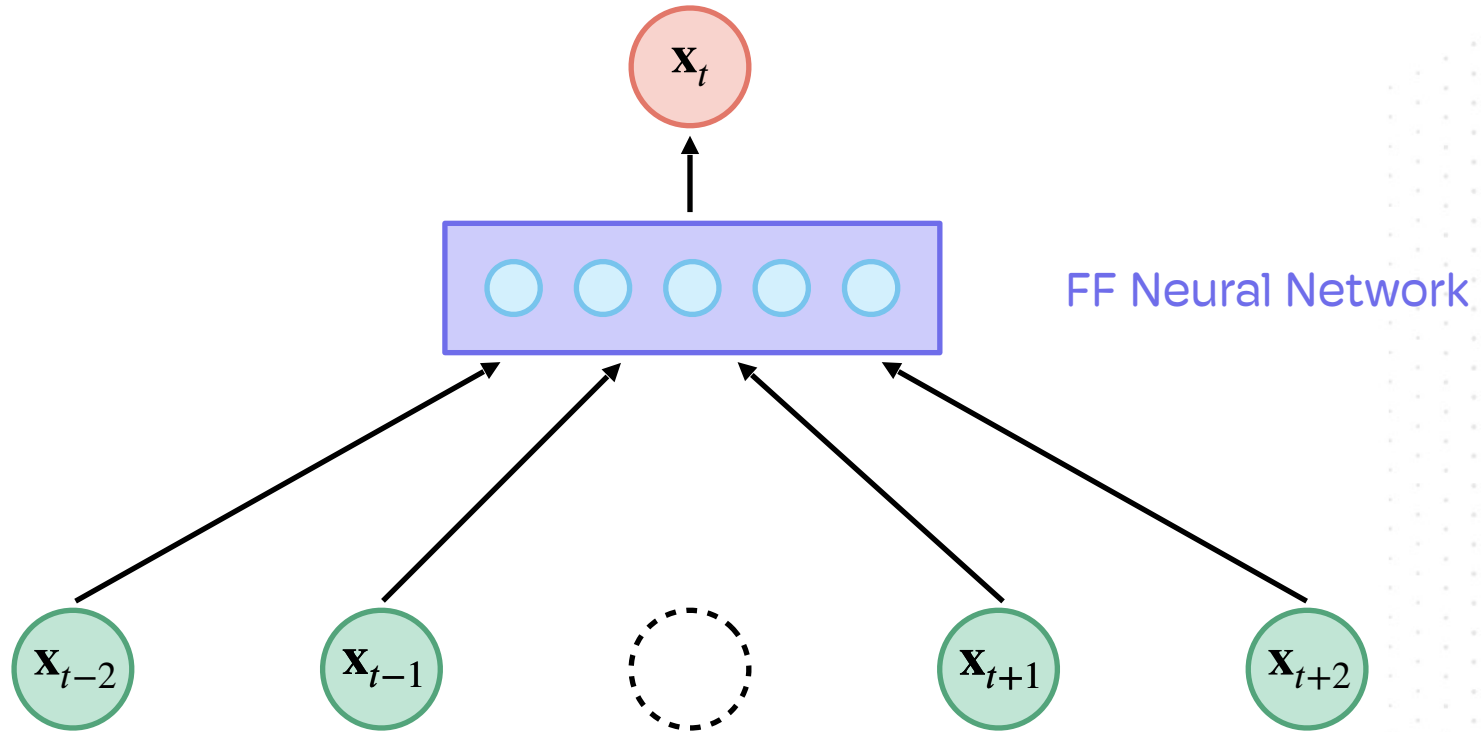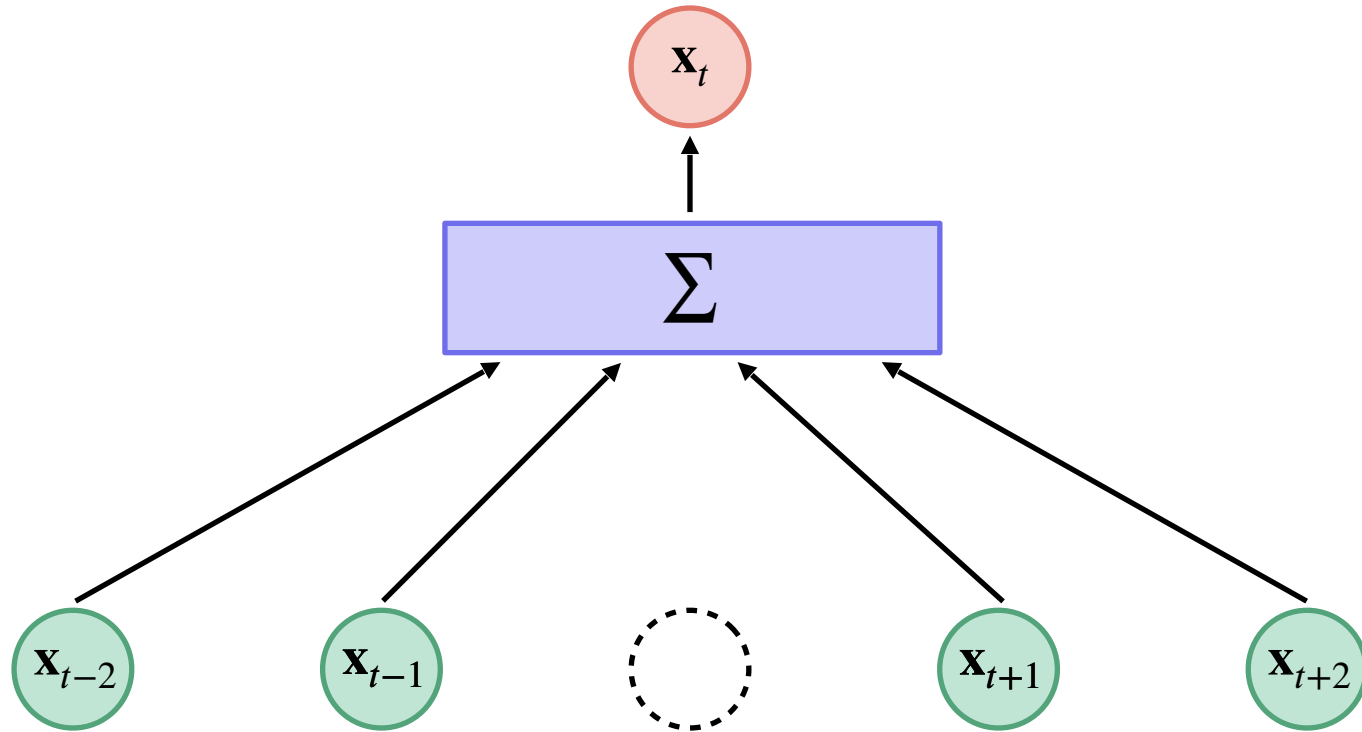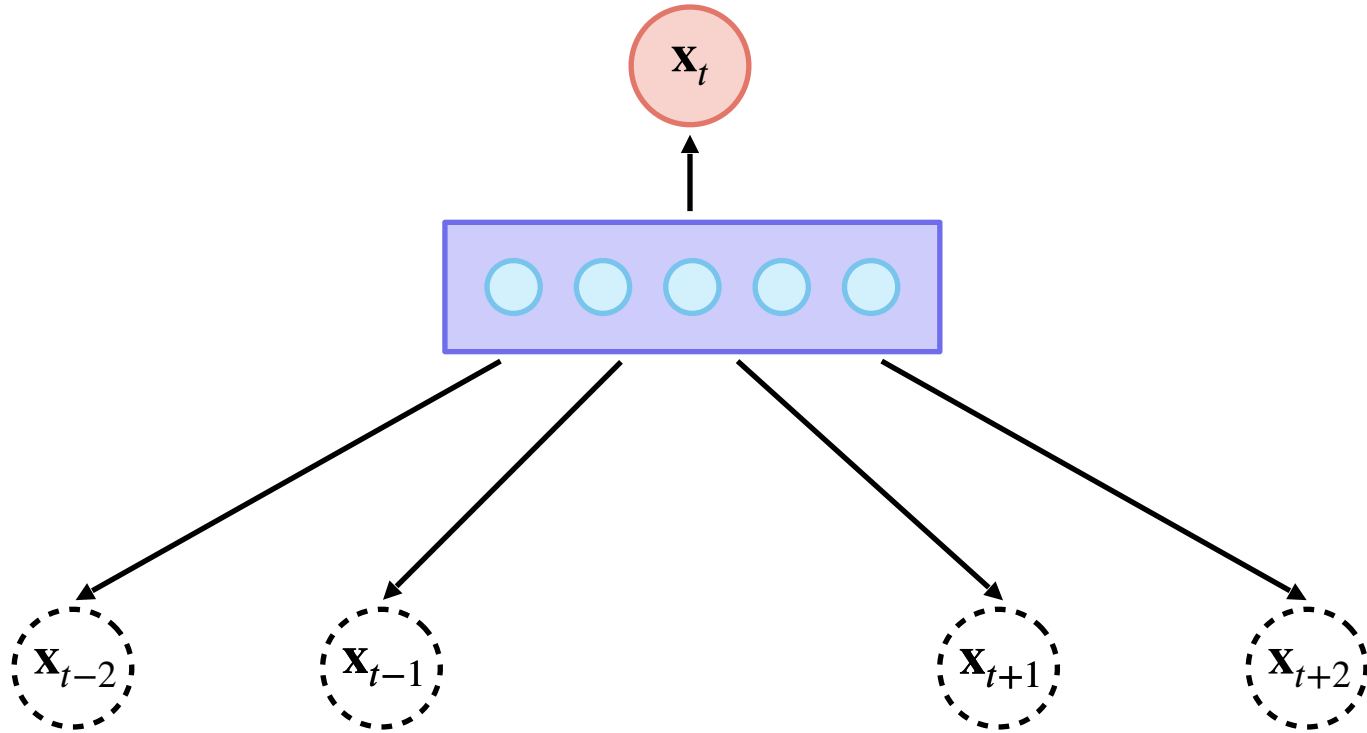
# Natural Language Manifold

# Natural Language Manifold

# Natural Language Manifold

# word2vec



FF Neural Network

# Skip-gram

# Continuous bag of words (CBOW)

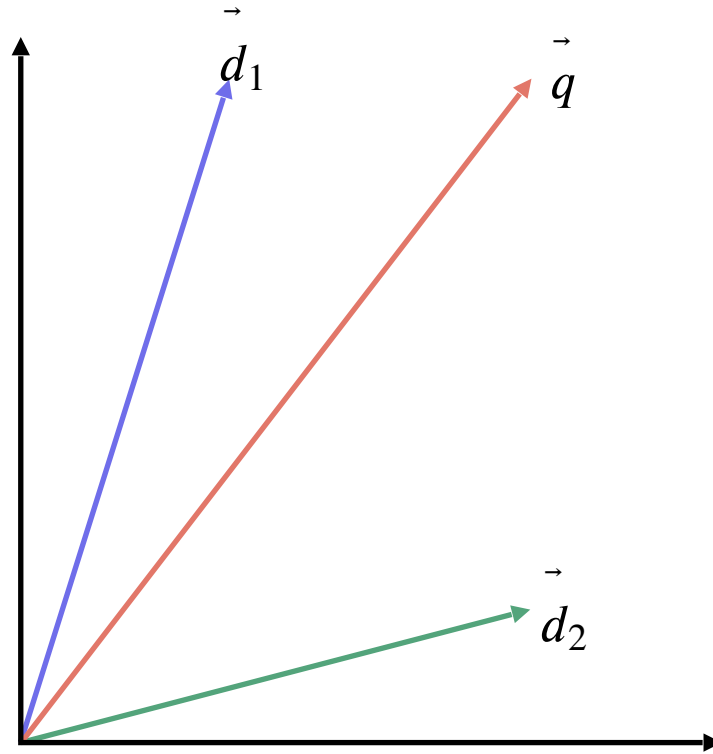# Embedding Visualizer
# Screen Sharing

# 5.9

## The Vector Space Model

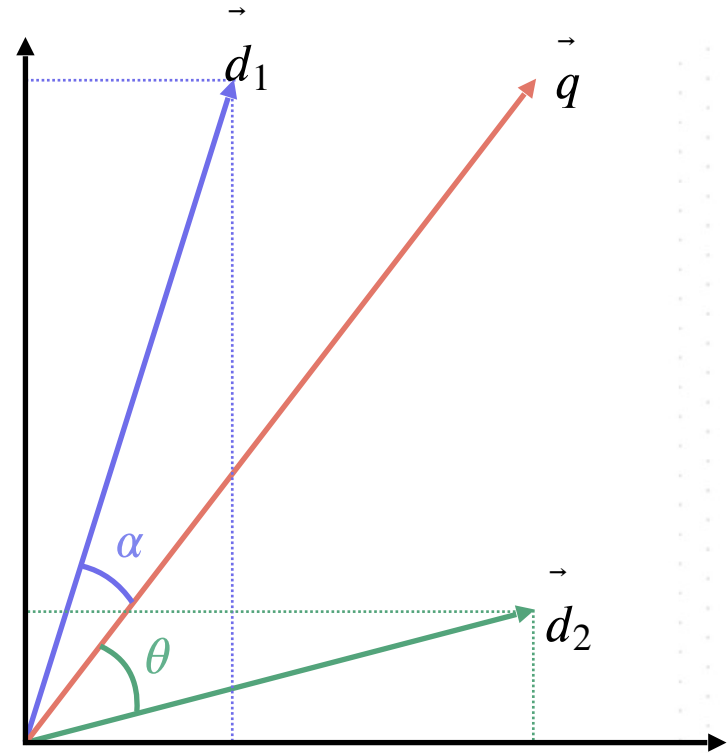# Unsupervised Learning

- No need for labels

- Discovers latent features (hidden patterns in data)

- Often exploratory in nature

- Since there is no "gold standard" often difficult to validate model (especially with stochastic algorithms)
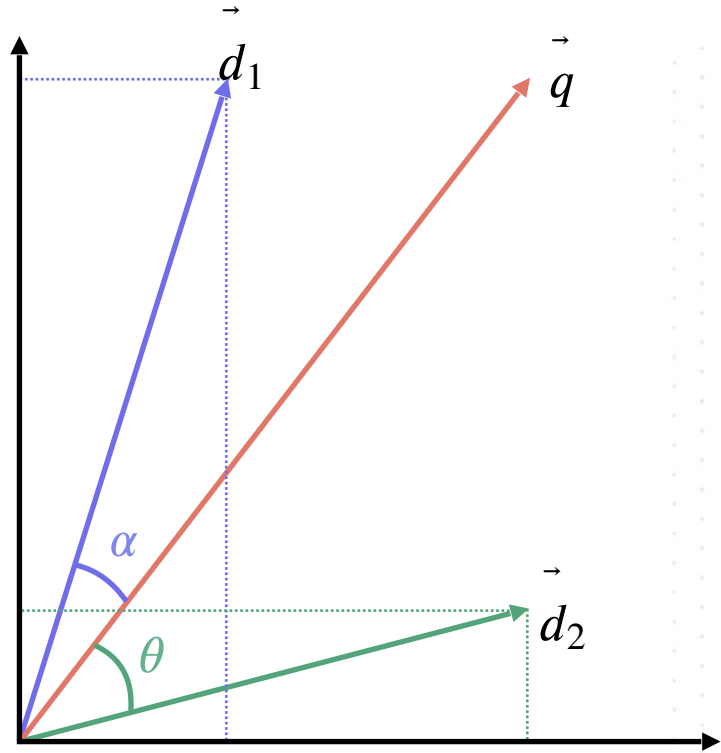
# Vector Space Model

# Vector Space Model

Similarity is a measure of "distance"

$$\cos\theta = \frac{\vec{d_2} \cdot \vec{q}}{\| \vec{d_2} \| \; \| \vec{q} \|}$$

**K-means ≈ PCA ≈ LDA ≈ SVD ≈ NMF**

Pearson

# It's All (Almost) the Same

Huang, Heng, et al. "Simultaneous tensor subspace selection and clustering: the equivalence of high order svd and k-means clustering." *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2008.

Ding, Chris, Xiaofeng He, and Horst D. Simon. "On the equivalence of nonnegative matrix factorization and spectral clustering." *Proceedings of the 2005 SIAM International Conference on Data Mining*. Society for Industrial and Applied Mathematics, 2005.

Ding, Chris, and Xiaofeng He. "K-means clustering via principal component analysis." *Proceedings of the Twenty-first International Conference on Machine learning*. 2004.

Corrochano, Eduardo Bayro, et al. "Eigenproblems in pattern recognition." *Handbook of Geometric Computing: Applications in Pattern Recognition, Computer Vision, Neuralcomputing, and Robotics* (2005): 129-167.

Pearson

# 5.10

## Embedding Sequences with Transformers

# Live Coding

# 5.11

## Computing the Similarity Between Embeddings

# Live Coding

# 5.12

## Semantic Search with Embeddings

# Live Coding

# 5.13

**Contrastive Embeddings with Sentence Transformers**

# Live Coding

Pearson