

A decorative network diagram in the top-left corner of the slide. It features a complex web of interconnected nodes and edges. The nodes are represented by small circles, some of which are solid blue, some are solid grey, and some are hollow with a blue outline. The edges are thin grey lines connecting the nodes. The overall pattern is dense and organic, resembling a molecular structure or a data network.

# **BST 261: Data Science II**

## **Lecture 16**

### **Course Review**

**Heather Mattie**  
**Harvard T.H. Chan School of Public Health**  
**Spring 2 2021**

A decorative network diagram in the bottom-right corner of the slide. It features a complex web of interconnected nodes and edges. The nodes are represented by small circles, some of which are solid blue, some are solid grey, and some are hollow with a blue outline. The edges are thin grey lines connecting the nodes. The overall pattern is dense and organic, resembling a molecular structure or a data network.

# Recipe of the Day!

Aperol Spritz Cocktail

Non-alcoholic Aperol Spritz



The background of the slide features a complex, light gray network pattern. It consists of numerous small circles, some of which are solid gray and others are hollow with a gray outline. These circles are interconnected by a web of thin, light gray lines, creating a dense, interconnected mesh that resembles a molecular structure or a data network. The overall tone is light and technical.

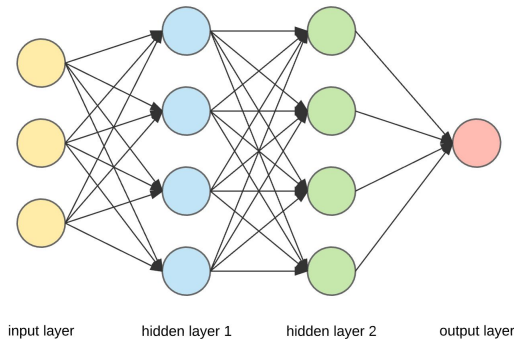
# Course Review

# What have we learned?

Quite a bit!

◎ Multilayer perceptrons → CNNs → RNNs → Advanced architectures

- ◎ Network architecture
- Hidden units
  - Layers
  - Activation function
  - Loss functions
  - Optimization algorithms
  - Batch size



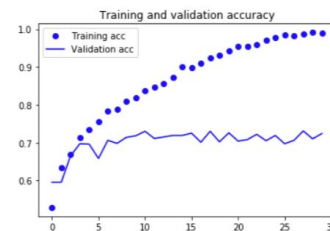
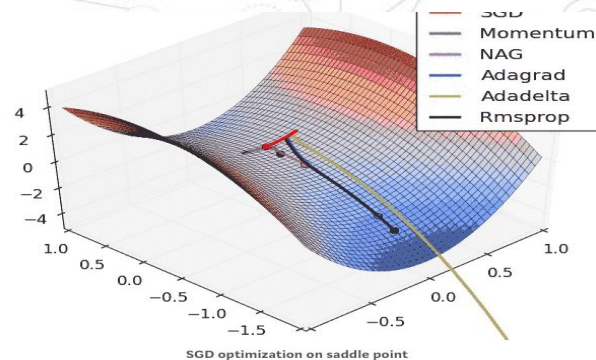
# What have we learned?

## How networks learn

- Gradient descent/ascent
- Backpropagation
- Forward pass and backward pass
- Visualizing filters
- Dense layer vs convolution layer vs RNN/LSTM/GRU layer

## Model performance

- Underfitting vs overfitting
- Regularization techniques
  - Dropout, L2 and L1 norms, network size
- Bias/variance tradeoff

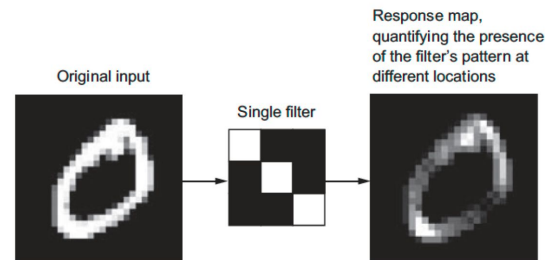
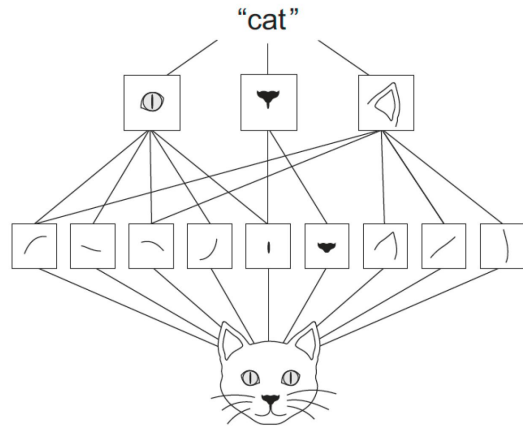


# What have we learned?



## CNNs

- Padding
- Pooling
- Strides
- Filters
- Translation invariance
- Hierarchical learning
- Lower layer representations vs higher layer representations
- Data format (3D vs 4D tensors)
- Object detection and localization
- Face recognition
- 1D CNN for sequential data
- Landmark detection
- Data augmentation
- Neural style transfer



# What have we learned?

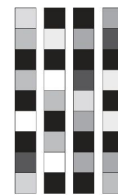


## RNNs

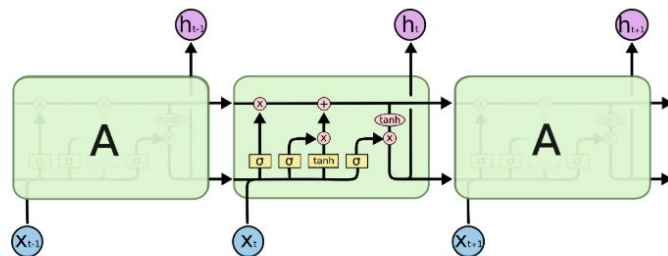
- Different types of sequential data
- How RNNs preserve order in this type of data
- SimpleRNN vs LSTM vs GRU layers
- Tokens and tokenization
- One-hot encoding and hashing
- Word embeddings
- Word2Vec and Glove
- Time series data
- Recurrent dropout
- Text generation
- Bidirectional recurrent layers



One-hot word vectors:  
- Sparse  
- High-dimensional  
- Hardcoded



Word embeddings:  
- Dense  
- Lower-dimensional  
- Learned from data



The repeating module in an LSTM contains four interacting layers.

# What have we learned?

## ◎ Advanced network architectures

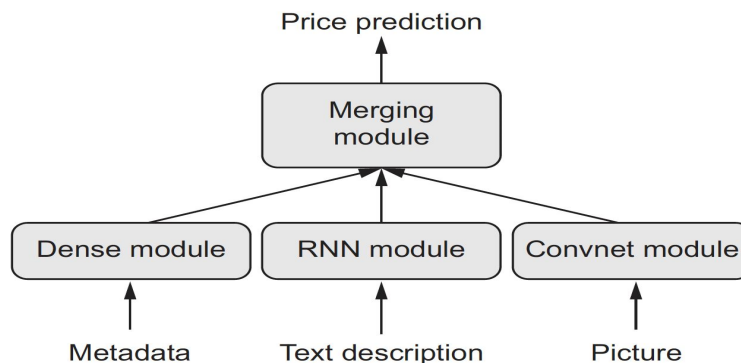
- One-to-many (multi-output/multi-head models)
- Many-to-many
- Many-to-one (multi-modal models)
- Directed acyclic graphs

## ◎ Advanced architecture patterns

- Batch normalization
- Hyperparameter optimization
- Model ensembling

## ◎ Implementation in Keras

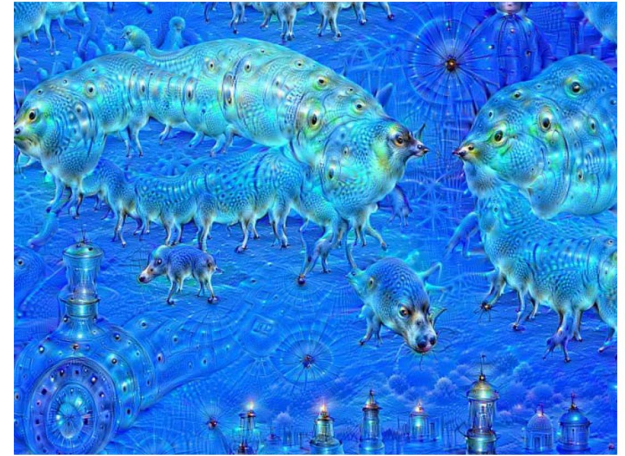
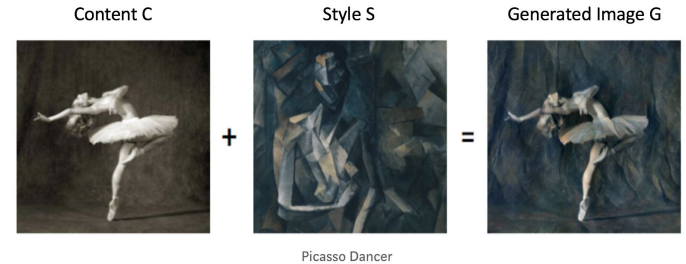
- Tensor (data) manipulation
- Sequential model
- MLP, CNN, RNN
- Using GCP
- Functional API





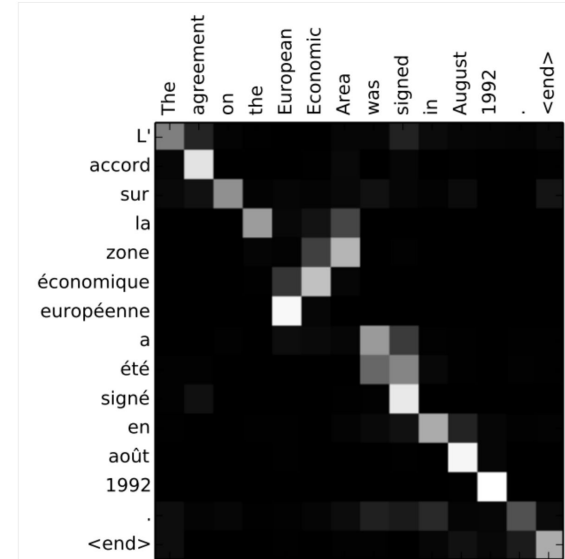
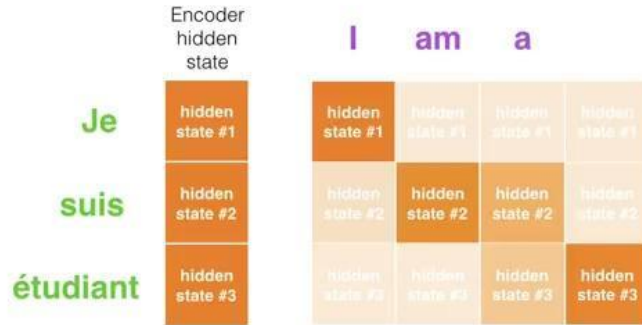
# What have we learned?

- ◎ Advanced topics
  - Variational autoencoders (VAEs)
  - Generative adversarial networks (GANs)
  - Reinforcement learning (RL)
  - DeepDream
  - Neural style transfer
  - Text generation



# Attention

- Scoring is done by the decoder at each time step
  - For each output word, scoring maps important / relevant words from the input sequence - higher weight means more relevance
  - This helps with the accuracy of the output prediction



You can see how the model paid attention correctly when outputting "European Economic Area". In French, the order of these words is reversed ("européenne économique zone") as compared to English. Every other word in the sentence is in similar order.

# Transformers

- Solve all of the problems with classic RNNs
  - Allow for parallel computing
  - Use **attention**
    - Helps with loss of information problem
- [Attention is all you need paper](#)
  - December 2017
  - Huge breakthrough in NLP

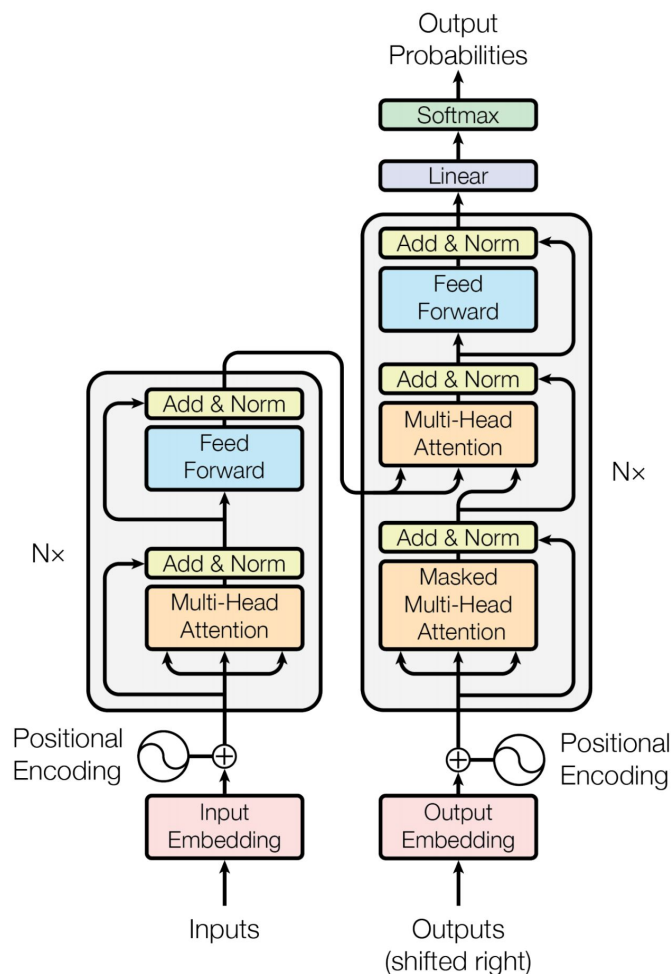


Figure 1: The Transformer - model architecture.

# Transformers

## State of the art models

- [GPT-2](#), [GPT-3](#)
- [BERT](#)
- [T5](#)

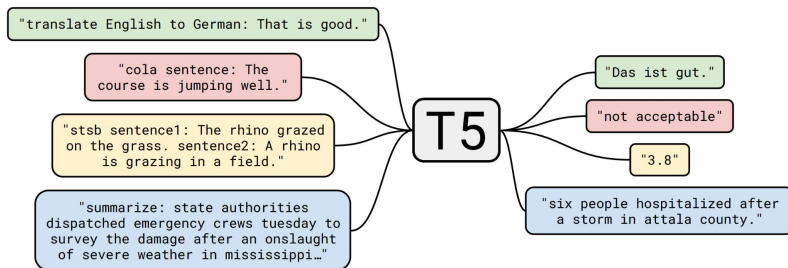
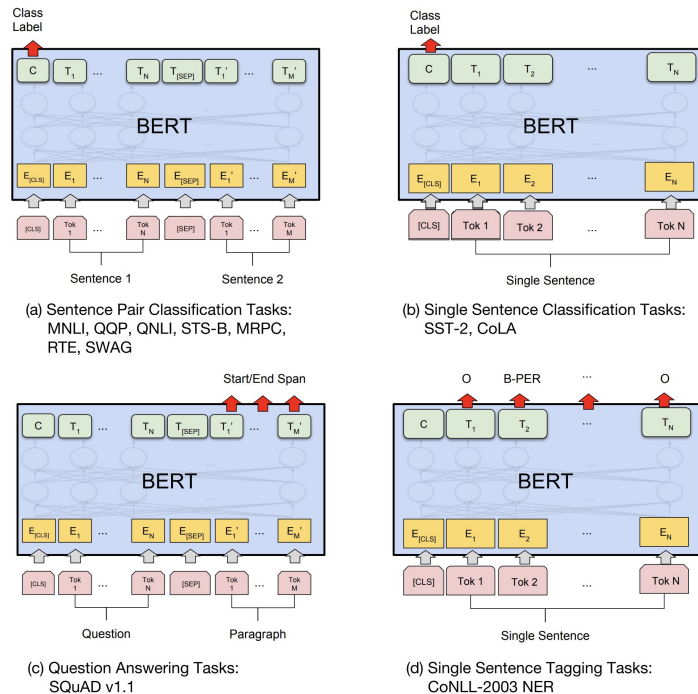
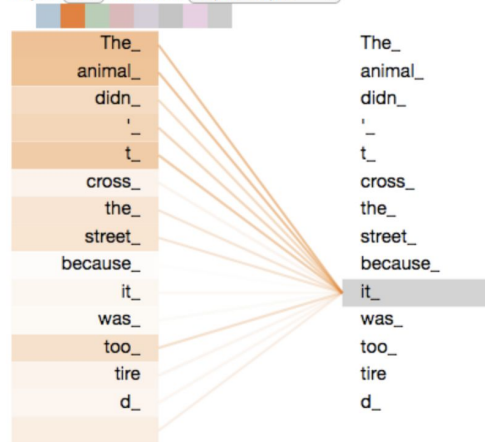


Figure 1: A diagram of our text-to-text framework. Every task we consider—including translation, question answering, and classification—is cast as feeding our model text as input and training it to generate some target text. This allows us to use the same model, loss function, hyperparameters, etc. across our diverse set of tasks. It also provides a standard testbed for the methods included in our empirical survey. “T5” refers to our model, which we dub the “Text-to-Text Transformer”.



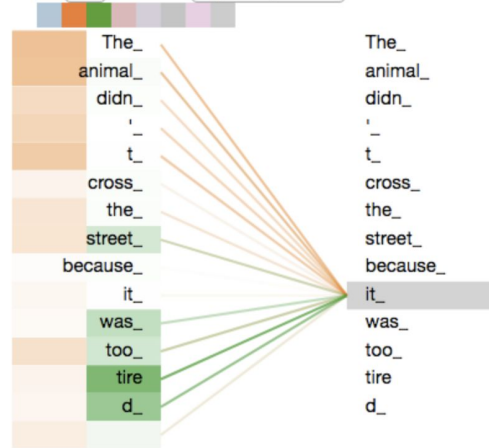
# Multi-Head Attention

Layer: 5 ▾ Attention: Input - Input ▾



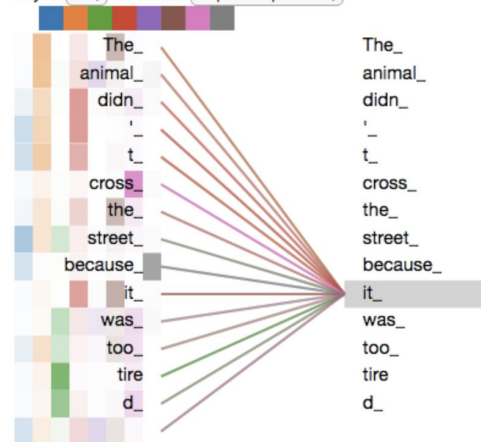
Self-attention

Layer: 5 ▾ Attention: Input - Input ▾



2-head Self-attention

Layer: 5 ▾ Attention: Input - Input ▾

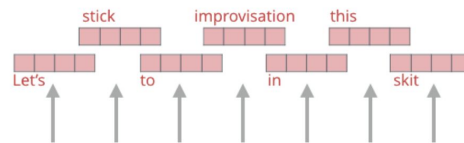


8-head Self-attention

# ELMo

- ◎ Recall: word embeddings are vector representations of words
  - We've used GloVe in class and in lab
  - Only 1 embedding for each word (fixed embeddings)
- ◎ Problem: the same word can have different meanings depending on the **context** in which it is used
- ◎ ELMo looks at an entire sequence before assigning each word it in an embedding
- ◎ Uses a bi-directional LSTM trained on a specific task

ELMo  
Embeddings



Words to embed



# Bidirectional Encoder Representations Transformer (BERT)

Has been described as the “marking the beginning of a new era in NLP”

- Has broken several records for language-based tasks
- Has been trained on massive datasets
- Was made available as a pre-trained network
- Now available in multiple sizes

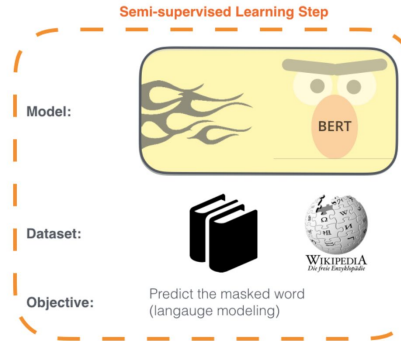
[Original paper](#)

[GitHub repository](#)

[Great post](#)

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



2 - **Supervised** training on a specific task with a labeled dataset.

