



# Midterm Information

- **Time: March 24th 5:45-7:15 PM**
- **Location (by section \*\*):**
  - Section 001 (Tuesday/Thursday 11-12:15PM): 6210 Social Sciences Bldg
  - Section002 (Tuesday/Thursday 2:30-3:45PM): B10 Ingraham Hall
  - Students with McBurney accommodations should have received an email with additional information.
  - Students who cannot take the exam on the specified time should contact their instructor if they have not done it yet.
- **Topics: Topics covered up to and including Week 9**
- **Exclusion List (questions regarding the following topics will NOT appear on the midterm):**
  - Logic (covered in sections 1 and 2)
  - SVM + Kernel Trick (covered in section 3)
- **Format:** MCQ
- **Cheat sheet:** a handwritten single piece of paper, front and back
- **Calculator:** optional, if it doesn't have an Internet connection
- **Bring:** your WISC ID, pencil (No 2 or softer), your 1-sheet notes.
- **Past exam questions:** on Canvas → Files → Past Exams

# Announcements

- **Homework 7 released and due Wednesday April 8 at 11:59PM.**
- **Additional optional review session on Friday March 20<sup>th</sup> at 11:00 AM on ZOOM (link provided on Canvas)**
- **TA Review session today at 5:30PM MH3610**
- **Additional Office hours: 3:45-5 PM Monday March 23rd**
- **Class roadmap and schedule:**

Machine Learning:  
Deep Learning III

ML/DL Summary

# Outline

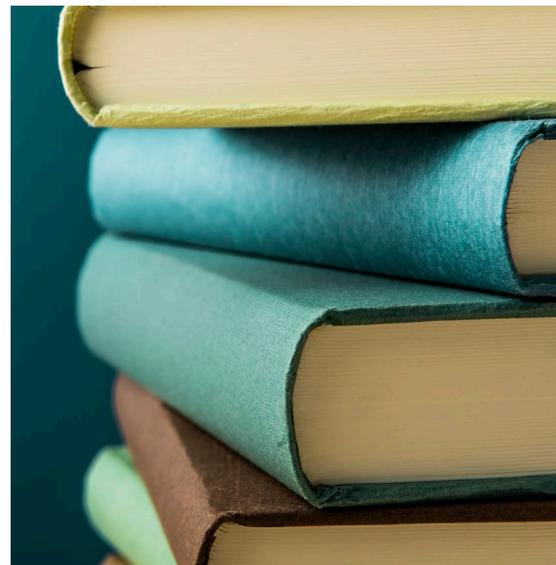
- Attention and Transformers
- Data Augmentation & Regularization
  - Expanding the dataset, avoiding overfitting
- More Signal From our Data
  - Graph-structured data, graph neural networks



# The Attention Mechanism

# From RNNS to Transformers

- RNNs process one word at a time — can't parallelize (Slow on GPUs.)
- Even LSTMs struggle with very long sequences
- All context gets squeezed into one hidden state vector



# From RNNS to Transformers

- Transformers allow parallelization and efficient context handling
- Example: “The cat the dog chased ran away.”
  - Who ran away?
  - Need to remember/attend to earlier words
- Goal: **decide which words matter most!**



# Word Representations & Context

- We use vectors to represent words (“embeddings”)

- Recall:

- One-hot representation

“dog”

[0 1 0 0 0 0 0 0 0]

- Dense embedding

- Vector captures **meaning** [0.13 0.87 - 0.23 0.46]



**Problem:** the meaning of a word depends on the words around it

# Word Representations & Context

“The monkey ate the banana. It was \_\_\_\_\_, wasn’t it?”

Could be:

- The monkey ate the banana. It was **ripe**, wasn’t it?
- The monkey ate the banana. It was **hungry**, wasn’t it?

Does “it” mean  
monkey or  
banana?

Meaning depends on past words (**context**)

The attention mechanism produces **contextual embeddings**.

# Attempt 1: Naïve Contextual Embedding

- Each token has a fixed embedding vector  $x_i$
- A crude attempt at contextual embedding: average over context
- Equal “attention” to every previous token

Tokens		Fixed Embeddings
1	the	[0.45 0.23]
2	monkey	[0.39 0.72]
3	ate	[0.83 0.61]
4	the	[0.45 0.23]
5	banana	[0.25 0.18]
6	it	[0.63 0.41]
7	was	[0.63 0.41]
8	ripe	[0.70 0.67]
9	wasn't	[0.14 0.61]
10	it	[0.63 0.41]
		<hr/>
		[4.98 4.93]

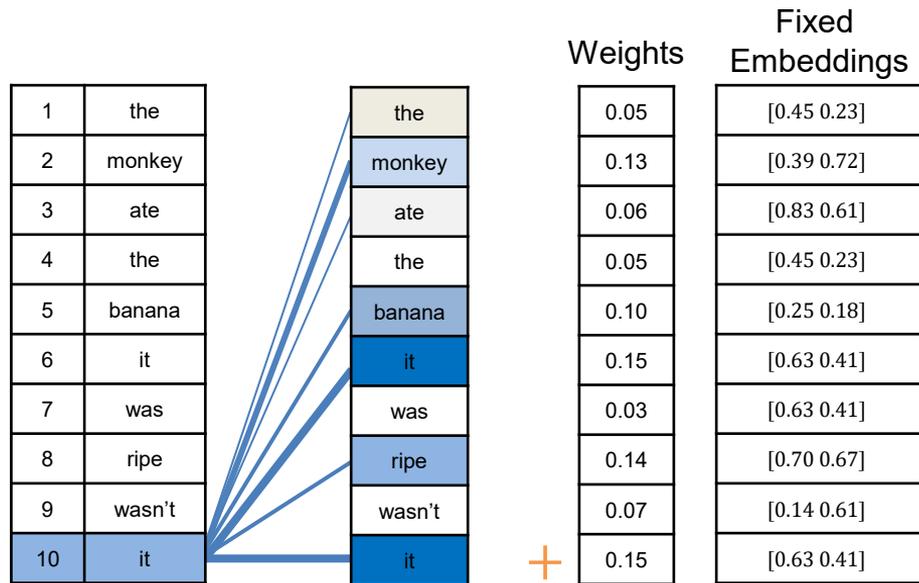
In math:  
for the  $i$ -th token

$$c_i = \frac{1}{i} \sum_{j=1}^i x_j$$

Contextual embedding for “it”: [0.498 0.493]

# Attempt 2: Assigning Weights

- Humans focus selectively
  - machines can too
- We can assign weights based on relevance
  - Idea: weight similar words highly
  - If  $\langle x_i, x_j \rangle$  large, assign large weight
- Then take weighted sum



In math: for  $i$ -th token

$$r_{ij} = \frac{1}{\sqrt{d}} \langle x_i, x_j \rangle$$
$$p_{i,:} = \text{softmax}(r_{i,:})$$

$$c_i = \sum_{j=1}^i p_{ij} \cdot x_j$$

# Final Attempt: The Attention Mechanism

Previous attempt:

- Fixed embeddings in three locations.


$$r_{ij} = \frac{1}{\sqrt{d}} \cdot \langle x_i, x_j \rangle$$

$$p_{i,:} = \text{softmax}(r_{i,:})$$


$$c_i = \sum_{j=1}^i p_{ij} \cdot x_j$$

In the **attention mechanism**:

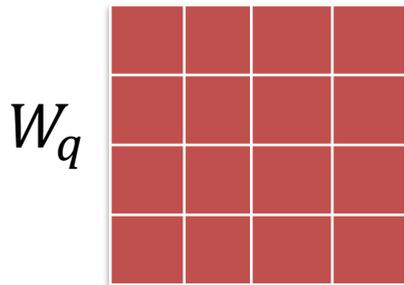
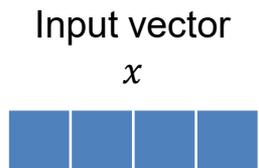
- Each token is associated with **three** vectors
- **Query**:  $q_i$ , the one attended from
- **Key**:  $k_j$ , the one attended to
- **Value**:  $v_j$ , the context being generated


$$r_{ij} = \frac{1}{\sqrt{d}} \cdot \langle q_i, k_j \rangle$$

$$p_{i,:} = \text{softmax}(r_{i,:})$$

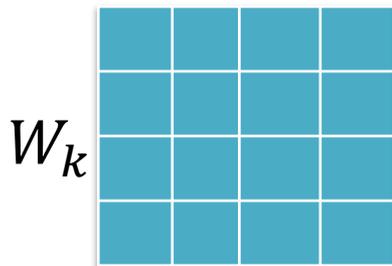

$$c_i = \sum_{j=1}^i p_{ij} \cdot v_j$$

# Query, Key, and Value Matrices



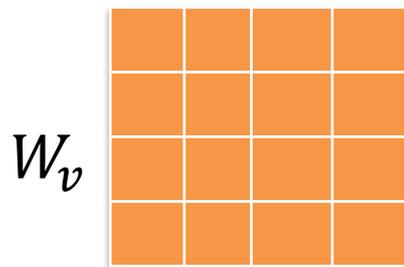
$$q = W_q x$$


Query



$$k = W_k x$$


Key

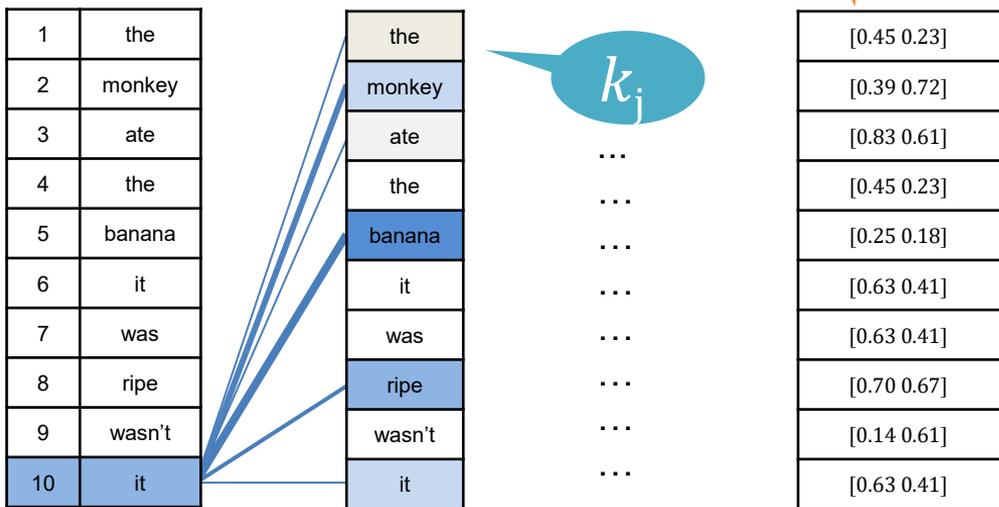


$$v = W_v x$$


Value

# The Attention Mechanism

Each token attends to all other tokens in the same sequence



$$r_{ij} = \frac{\langle q_i, k_j \rangle}{\sqrt{d}}$$



$$p_{i,:} = \text{softmax}(r_{i,:})$$



$$c_i = \sum_j p_{ij} \cdot v_j$$

# Notation for Attention

Queries, keys and values are written as matrices  $\mathbf{Q}, \mathbf{K}, \mathbf{V}$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

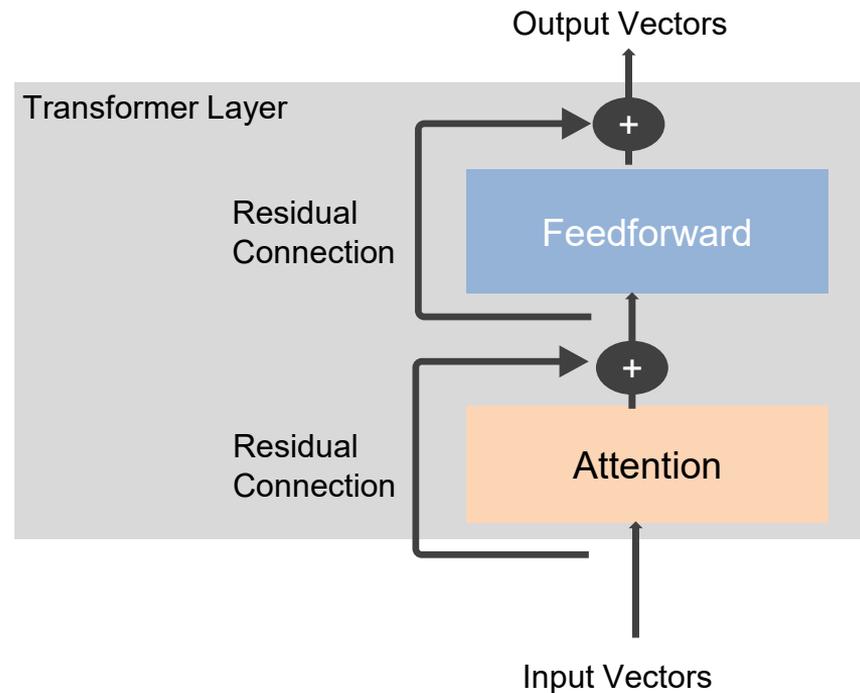


# The Transformer Architecture

# From Attention to Transformer

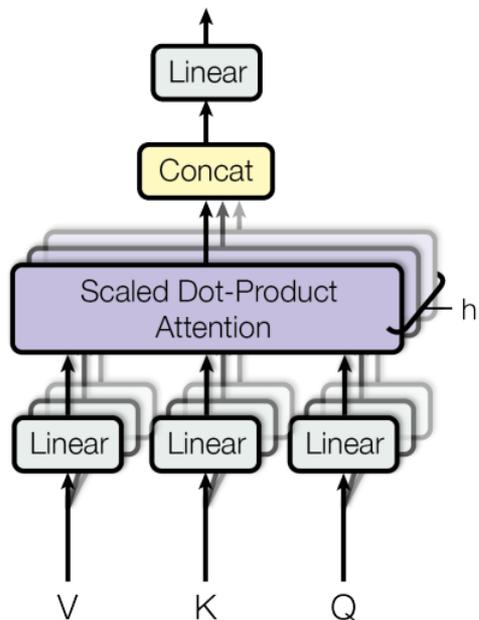
A single layer transformer consists of:

- Attention Mechanism
- Feed-Forward Network
- Residual Connections



# Multi-Head Attention

- Outputs combined for richer representations
- Multiple heads learn different relationships (syntax, meaning, position)



# Positional Encoding

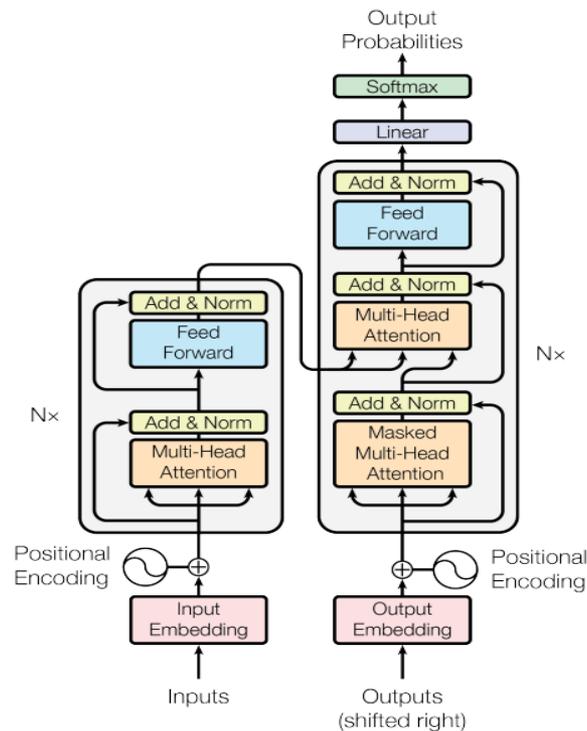
- Transformers have no recurrence — order must be added explicitly
- **Positional Encoding:** Information about the relative or absolute position of the tokens in the sequence
- Added to the input embeddings

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

position      dimension index      dimension

# Transformer Architecture

- Encoder–Decoder structure
- **Encoder:** maps an input sequence to a sequence of continuous representations  $z$ .
  - Useful for classification
- **Decoder:** Given  $z$ , the decoder generates an output sequence of symbols one element at a time.
  - Useful for generation

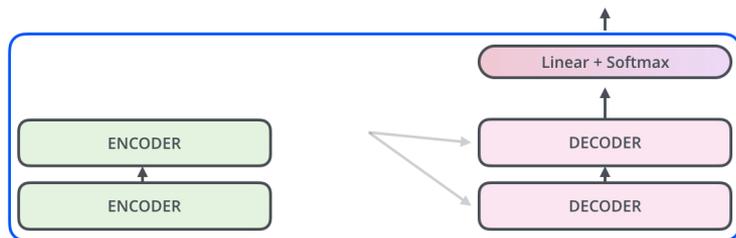


# Decoder

- **Masked multi-head attention:** each word attends to the words before it
- A **second attention** module that **attends the output of the encoder**

Decoding time step: 1 2 3 4 5 6

OUTPUT



EMBEDDING WITH TIME SIGNAL

	■ ■ ■ ■	■ ■ ■ ■	■ ■ ■ ■
--	---------	---------	---------

EMBEDDINGS

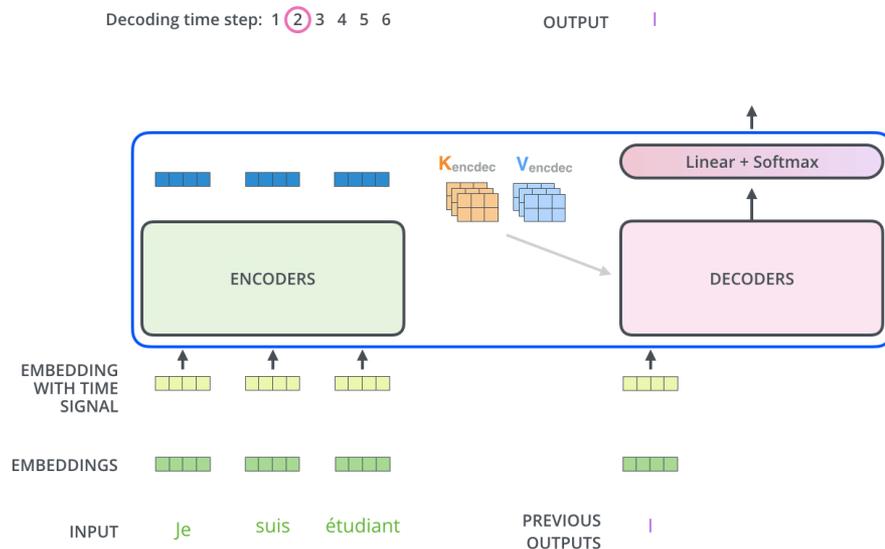
	■ ■ ■ ■	■ ■ ■ ■	■ ■ ■ ■
--	---------	---------	---------

INPUT

	Je	suis	étudiant
--	----	------	----------

# Decoder

- **Masked multi-head attention:** each word attends to the words before it
- A **second attention** module that **attends the output of the encoder**



## Q3.1 Quiz Break

What is the primary function of the self-attention mechanism?

- A) To track the position of each word in the sequence.
- B) To process the input sequence strictly from left to right.
- C) To weigh the importance and relationship of all words in a sequence relative to each other.
- D) To reduce the size of the model by using fewer parameters.

## Q3.1 Quiz Break

What is the primary function of the self-attention mechanism?

- A) To track the position of each word in the sequence.
- B) To process the input sequence strictly from left to right.
- C) **To weigh the importance and relationship of all words in a sequence relative to each other.**
- D) To reduce the size of the model by using fewer parameters.

## Q3.2 Quiz Break

In the context of the Transformer model, what role do "Encoders" and "Decoders" play?

- A) Encoders are used for text generation, and Decoders are used for text classification.
- B) Encoders process the input sequence, and Decoders generate the output sequence.
- C) Both Encoders and Decoders are only used for understanding the input sequence.
- D) Both Encoders and Decoders map an input sequence to a sequence of continuous representations.

## Q3.2 Quiz Break

In the context of the Transformer model, what role do "Encoders" and "Decoders" play?

- A) Encoders are used for text generation, and Decoders are used for text classification.
- B) Encoders process the input sequence, and Decoders generate the output sequence.**
- C) Both Encoders and Decoders are only used for understanding the input sequence.
- D) Both Encoders and Decoders map an input sequence to a sequence of continuous representations.

# Applications

- Language Models: GPT, BERT, T5
- Vision: ViT (Vision Transformer)
- Multimodal: CLIP, DALL·E, GPT-4v
- Scientific: AlphaFold, time-series modeling, robotics

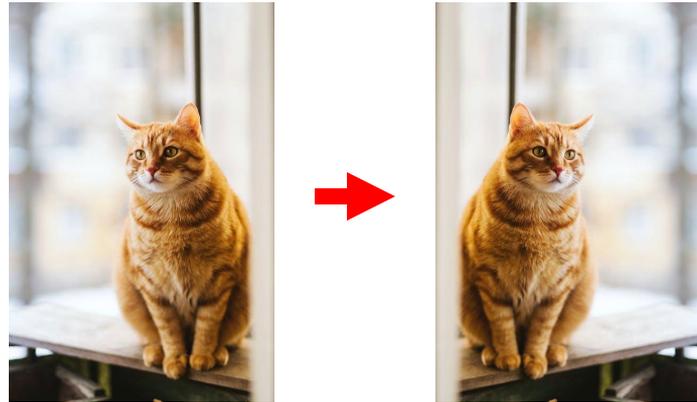


# Data Concerns

# Data Concerns

What if we don't have a lot of data?

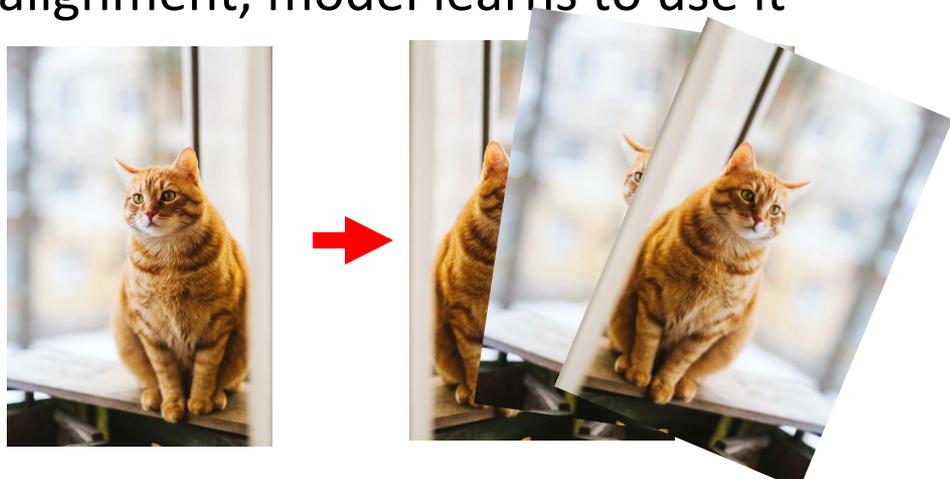
- We risk overfitting
- Avoiding overfitting: **regularization** methods
- Data augmentation: a classic way to regularize



# Data Augmentation

Augmentation: transform + add new samples to the training set

- Transformations: based on domain
- Idea: build **invariances** into the model
  - **Ex:** if all images have same alignment, model learns to use it
- Keep the label the same!



# Transformations

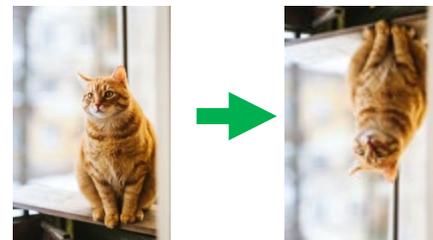
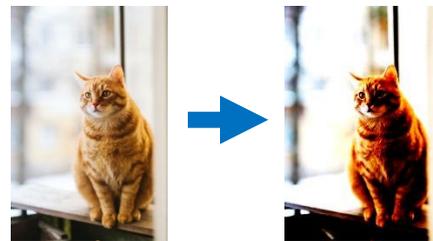
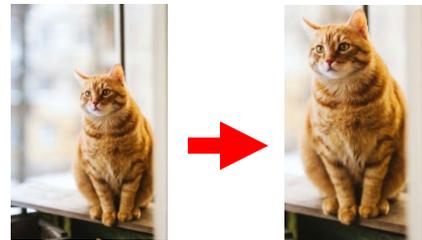
## Examples of transformations for images

- **Crop** (and zoom)
- **Color** (change contrast/brightness)
- **Rotations+** (translate, stretch, shear, etc)

Many more possibilities. Combine as well!

Q: how to deal with this at **test time**?

- A: transform, test, average



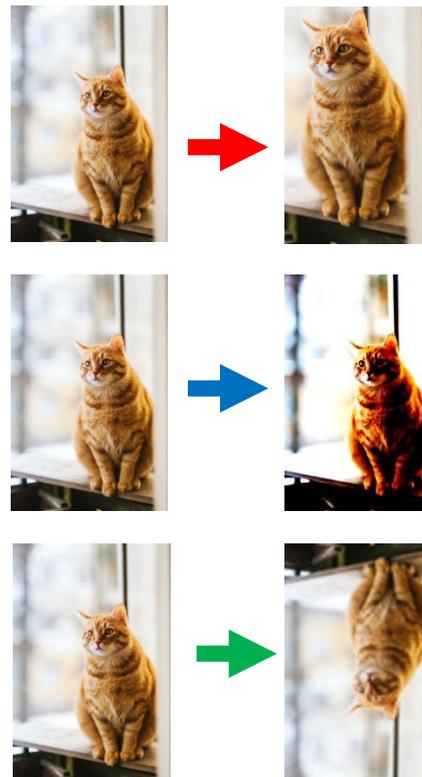
# Combining & Automating Transformations

One way to automate the process:

- Apply every transformation and combinations
- **Downside:** most don't help...

Want a good policy, ie, → → → → →

- Active area of research: search for good policies
  1. **Ratner et al:** "Learning to Compose Domain-Specific Transformations for Data Augmentation"
  2. **Cubuk et al:** "AutoAugment: Learning Augmentation Strategies from Data"



# Other Domains

Not just for image data. For example, on text:

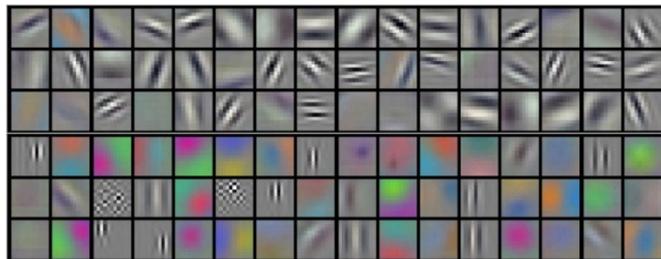
- Substitution
  - E.g., “It is a **great** day” → “It is a **wonderful** day”
  - Use a thesaurus for particular words
  - Or, use a model. Pre-trained word embeddings, language models
- Back-translation
  - “Given the low budget and production limitations, this movie is very good.”  
→ “There are few budget items and production limitations to make this film a really good one”

# Importance of Augmentation

Data augmentation is critical for top performance!

- You should use it!
- **AlexNet**: used (many papers re-used as well)
  - Random crops, rotations, flips.

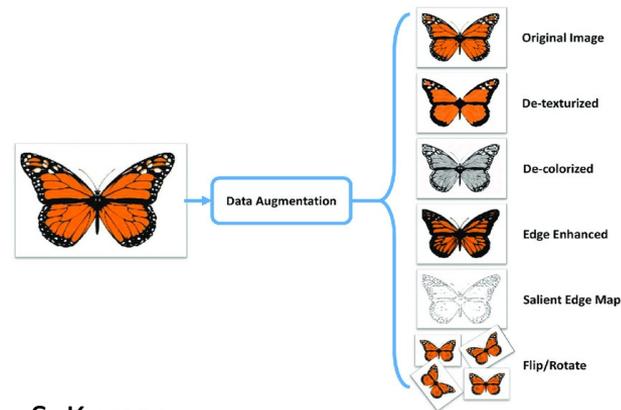
Krizhevsky et al: “ImageNet Classification with Deep Convolutional Neural Networks”



# Other Forms of Regularization

## Regularization has many interpretations

- **Goodfellow:** “any modification... to a learning algorithm that is intended to reduce its generalization error but not its training error.”
- A way of adding knowledge / side information to model
- Enforcing parsimony/simplicity



# Other Forms of Regularization

## Classic regularizations

1. Modify loss functions

**Ex:** regularized least squares LR

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n (\theta_0 + x_i^T \theta - y_i)^2 + \lambda \|\theta\|_2^2$$

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i) + \lambda R(f_{\theta})$$

Standard  
loss

Regularization  
parameter

Regularizer

2. Modify architecture/training/data

a) Dropout, batch normalization, augmentation

# Break & Quiz

**Q 2.1:** If we apply data augmentation blindly, we might

(i) Change the label of the data point

(ii) Produce a useless training point

- A. (i) but not (ii)
- B. (ii) but not (i)
- C. Neither
- D. Both

# Break & Quiz

**Q 2.1:** If we apply data augmentation blindly, we might

(i) Change the label of the data point

(ii) Produce a useless training point

A. (i) but not (ii)

B. (ii) but not (i)

C. Neither

**D. Both**

# Break & Quiz

**Q 2.1:** If we apply data augmentation blindly, we might

(i) Change the label of the data point

(ii) Produce a useless training point

A. (i) but not (ii) (Can do (ii): imagine turning up the contrast till the image is completely black and is unusable).

B. (ii) but not (i) (Can change label: rotate a 6 into a 9).

C. Neither (Can do either).

**D. Both**

# Break & Quiz

**Q 2.2:** Which of the following is false about regularization?

A. Data augmentation can be applied for text.

B. Regularization always introduces complexity to a model to avoid overfitting.

C. We can modify the loss function by adding a regularization term.

D. The transformations for data augmentation depend on the domain.

# Break & Quiz

**Q 2.2:** Which if the following is **false** about regularization?

A. Data augmentation can be applied for text.

**B. Regularization always introduces complexity to a model to avoid overfitting.**

C. We can modify the loss function by adding a regularization term.

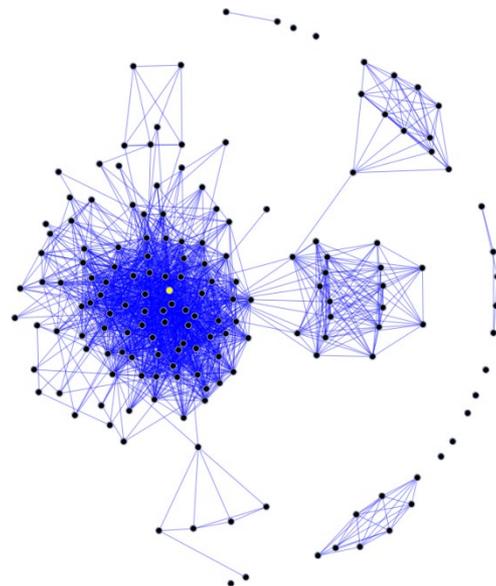
D. The transformations for data augmentation depend on the domain.



# Graph Neural Networks

# Relationships in Data

- We usually assume all data points are independent, “unrelated” in a sense
- Pretty common to have relationships between points  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ 
  - **Social networks**: individuals related by friendship
  - **Biology/chemistry**: bonds between compounds, molecules
  - **Citation networks**: Scientific papers cite each other

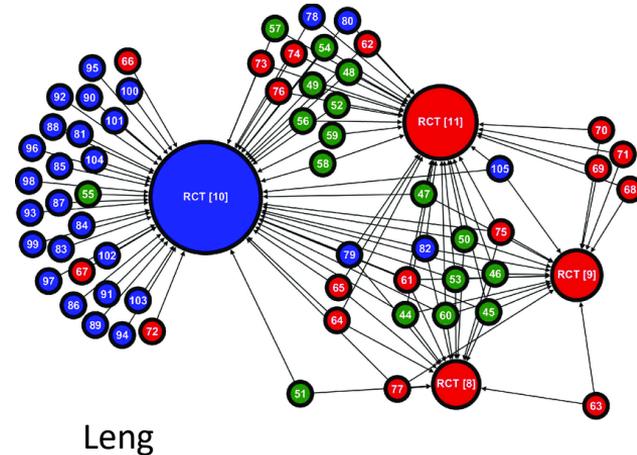


Wiki

# Signal from Relationships

Suppose we are classifying scientific papers

- **Features:** title, abstract, authors. **Labels:** math/science/eng.
- Could build a reasonable classifier with the above data
- **More signal** from relationships
  - Cite each other, more likely from the same field
  - Note: citations are not features; they're **links**
  - Need a new type of network to handle



# Graph Neural Networks

Have:  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n), G = (V, E)$

How should our new architecture look?

- Still want layers
  - linear transformation + non-linearity
- Now want to integrate neighbors
- Bottom: graph convolutional network

Hidden Layer Representation



$$H^{(\ell+1)} = \sigma(H^{(\ell)}W^{(\ell)})$$

Non-Linearity



Parameters



$$H^{(\ell+1)} = \sigma(A_G H^{(\ell)} W^{(\ell)})$$

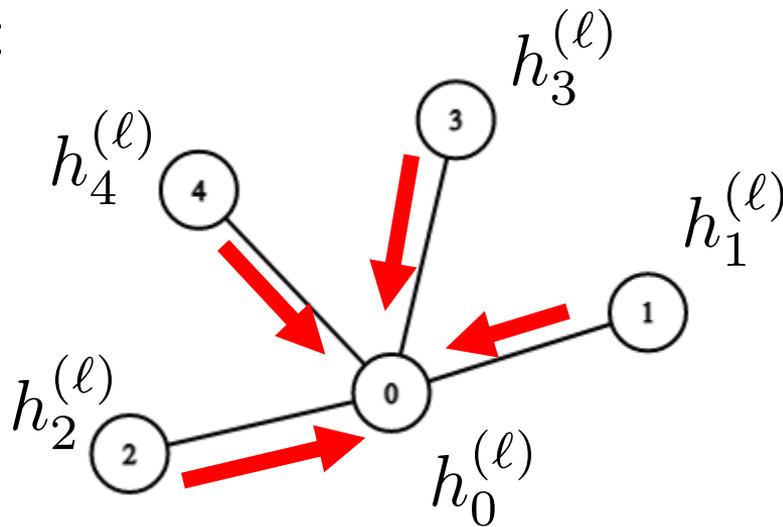


Graph Mixing

# Graph Convolutional Networks

Let's examine the GCN architecture in more detail

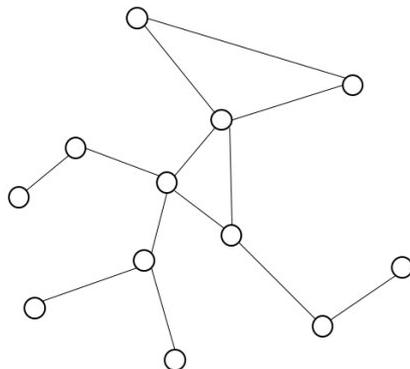
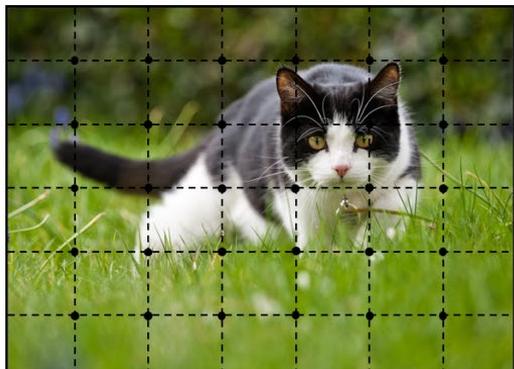
- Difference: “graph mixing” component
- At each layer, get representation at each node
- Combine node's representation with neighboring nodes
- **“Aggregate”** and **“Update”** rules



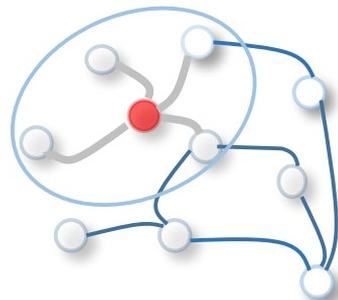
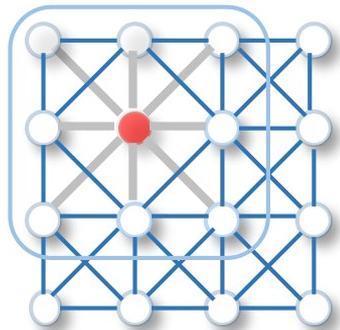
# Graph Convolutional Networks

Note the resemblance to CNNs:

- Pixels: arranged as a very regular graph
- Want: more general configurations (less regular)



Wu et al, A Comprehensive Survey on Graph Neural Networks



# Further Reading/Viewing

- Jurafsky & Martin, Chapter 8
  - [https://web.stanford.edu/~jurafsky/slp3/ed3book\\_aug25.pdf](https://web.stanford.edu/~jurafsky/slp3/ed3book_aug25.pdf)
- Russell & Norvig, Chapter 21.6 and 24
- Andrej Karpathy tutorial
  - <https://karpathy.ai/zero-to-hero.html>
- 3Blue1Brown:
  - <https://www.youtube.com/watch?v=eMlx5fFNoYc>
- The Illustrated Transformer
  - <https://jalammar.github.io/illustrated-transformer/>

