Cornell Bowers C·IS
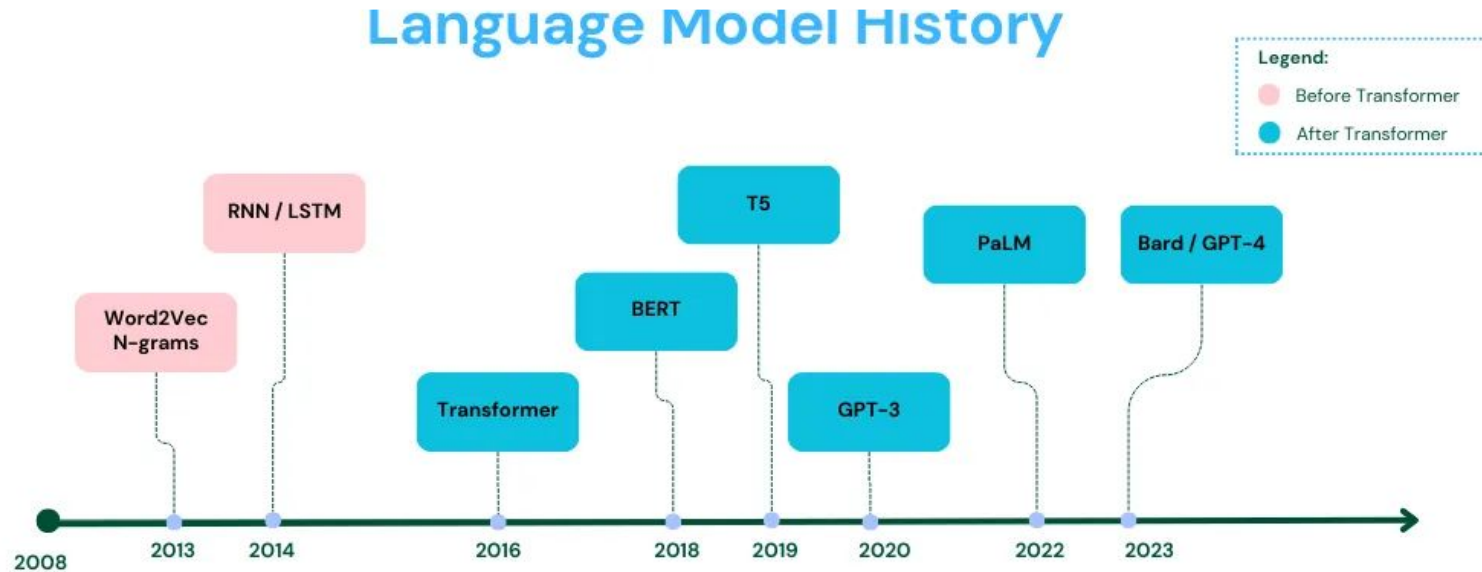College of Computing and Information Science

Transformers

# How to handle text data?

# Language Modelling History



**Language Model History**

Legend:
- Before Transformer
- After Transformer

RNN / LSTM

Word2Vec N-grams

Transformer

BERT

T5

GPT-3

PaLM

Bard / GPT-4

2008 — 2013 — 2014 — 2016 — 2018 — 2019 — 2020 — 2022 — 2023

# Language Modeling: predict the next word

**Assign probabilities to text.**

Given a sequence $(\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$, we want to **maximize** $P(\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T)$.

$$P(\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T) = P(\mathbf{x}_1)P(\mathbf{x}_2|\mathbf{x}_1)P(\mathbf{x}_3|\mathbf{x}_1, \mathbf{x}_2)P(\mathbf{x}_4|\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \ldots P(\mathbf{x}_T|\mathbf{x}_1, \mathbf{x}_2, \ldots \mathbf{x}_{T-1})$$

P(I like cats because they look cute) = P(I) P(like | I) P(cats | I like) P(as | I like cats)  P(they | I like cats because)

P(look | I like cats because they)  P(cute | I like cats because they look)

Predict the next word given current text!

# *n*-Gram Language Model

*n*-Gram: chunk of *n* consecutive words

**Count** the frequency of each *n*-grams and predict next word!

**Uni-gram:** "I" "like" "cats" "as" "they" "look" "cute"

**Bi-gram:** "I like" "like cats" "cats as" "as they" …

**Tri-gram:** "I like cats" "like cats as" "cats as they" …

**Assume** each word only depends on previous *n - 1* words.

$$P(\mathbf{x}_t|\mathbf{x}_1,\ldots,\mathbf{x}_{t-1}) = P(\mathbf{x}_t|\mathbf{x}_{t-n+1},\ldots,\mathbf{x}_{t-1})$$
$$= \frac{\text{count}(\mathbf{x}_{t-n+1},\ldots,\mathbf{x}_{t-1},\mathbf{x}_t)}{\text{count}(\mathbf{x}_{t-n+1},\ldots,\mathbf{x}_{t-1})}$$

In *bi-gram* LM
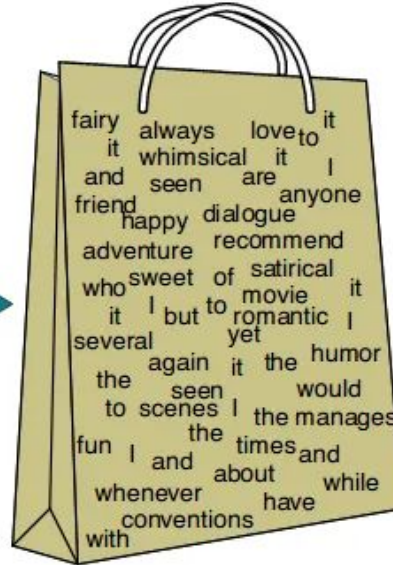
P(I like cats as they look cute) = P(I) P(like | I) P(cats | like) P(as | cats)  P(they | because)  P(look | they)  P(cute | look)

Discuss: Do you want to have a large n or a small n in a n-gram model?

# Bag of Words

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!
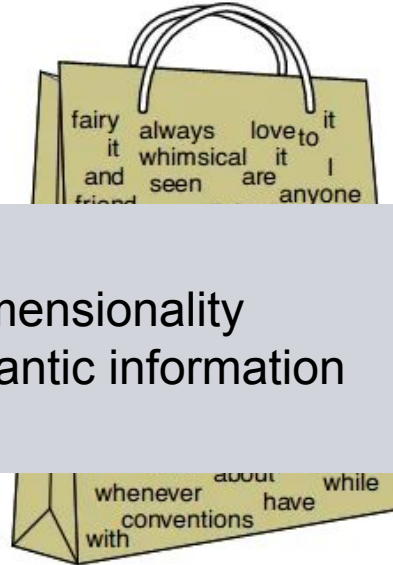
fairy always love to it
it whimsical it I
and seen are anyone
friend happy dialogue
adventure recommend
who sweet of satirical it
it I but to movie it
several yet romantic I
the again it the humor
to seen would
fun scenes I the manages
I the times and
and about while
whenever have
conventions
with

| it | 6 |
|----|---|
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| yet | 1 |
| would | 1 |
| whimsical | 1 |
| times | 1 |
| sweet | 1 |
| satirical | 1 |
| adventure | 1 |
| genre | 1 |
| fairy | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| … | … |

# Bag of Words

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to and romantic at the conven fairy tale geni recommend it anyone. I've s times, and I'm to see it again have a friend seen it yet!
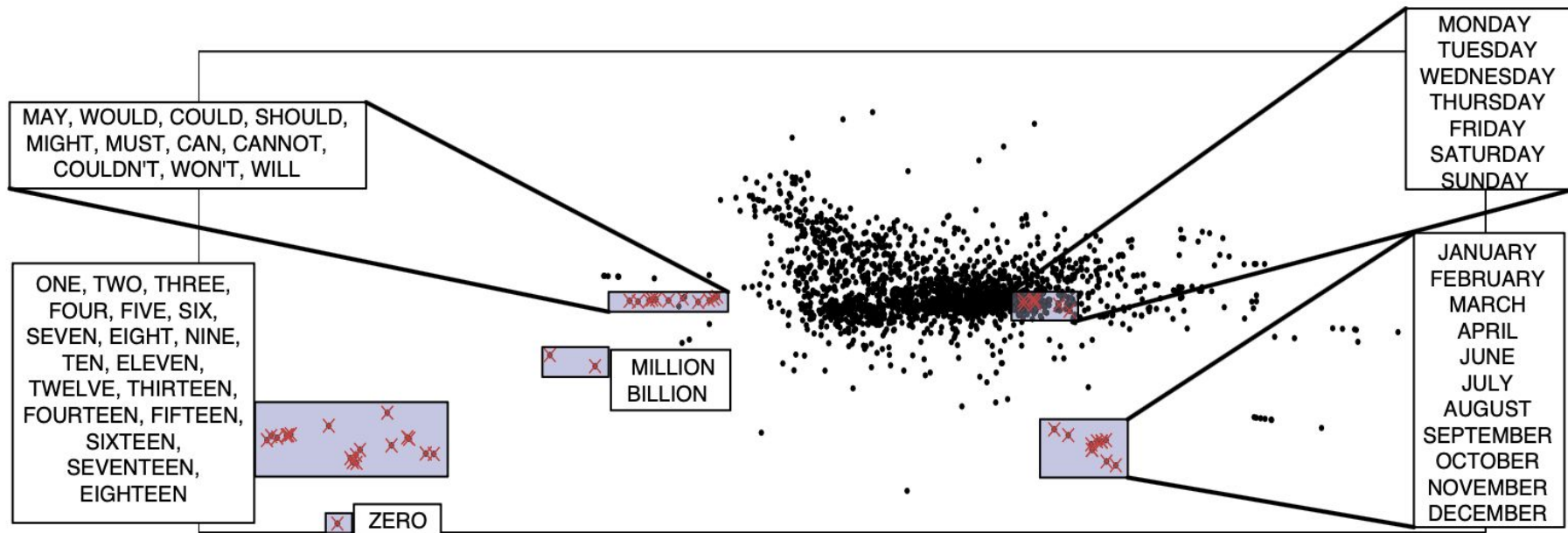
fairy always love to it
it whimsical it
and seen are I
friend anyone
whenever about while
conventions have
with

| it | 6 |
|---|---|
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| | 1 |
| | 1 |
| sical | 1 |
| | 1 |
| al | 1 |
| ture | 1 |
| | 1 |
| | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| ... | ... |

Drawbacks:
● High dimensionality
● No semantic information

# Word Embeddings



MONDAY
TUESDAY
WEDNESDAY
THURSDAY
FRIDAY
SATURDAY
SUNDAY

MAY, WOULD, COULD, SHOULD,
MIGHT, MUST, CAN, CANNOT,
COULDN'T, WON'T, WILL

JANUARY
FEBRUARY
MARCH
APRIL
JUNE
JULY
AUGUST
SEPTEMBER
OCTOBER
NOVEMBER
DECEMBER

ONE, TWO, THREE,
FOUR, FIVE, SIX,
SEVEN, EIGHT, NINE,
TEN, ELEVEN,
TWELVE, THIRTEEN,
FOURTEEN, FIFTEEN,
SIXTEEN,
SEVENTEEN,
EIGHTEEN

MILLION
BILLION

ZERO

Blitzer et al. 2004

# What does ong choi mean?

Suppose you see these sentences:

- Ong choi is delicious sautéed with garlic.
- Ong choi is superb over rice
- Ong choi leaves with salty sauces

And you've also seen these:

- ...spinach sautéed with garlic over rice
- Chard stems and leaves are delicious
- Collard greens and other salty leafy greens

Ong choy is **a leafy green vegetable with long, hollow stems and slender leaves**. It's also known as Chinese water spinach, Chinese water spinach, or hollow stem spinach.

**S** The Seasoned Wok

**Ong Choy (Water Spinach) Recipe with Fermented Bean...**

Nov 3, 2022 — What is Ong Choy, Rau Muong or Water Spinach? Ong choy i...

The Woks of Life

**Ong Choy with XO sauce - The Woks of Life**

May 2, 2017 — Ong Choy is a popular Chinese leafy green vegetable that's...

Onolicious Hawai'i

**Garlic and Fish Sau... - Onolicious Hawaii**

Jan 7, 2021 — What is O Hawaii everyone know...

# Word2Vec

- We want vectors for words so that the context of a word can suggest the vector of this word, and vice versa
- Idea: **Similar words appear in similar contexts**

**Efficient Estimation of Word Representations in Vector Space**

**Tomas Mikolov**
Google Inc., Mountain View, CA
tmikolov@google.com

**Kai Chen**
Google Inc., Mountain View, CA
kaichen@google.com

**Greg Corrado**
Google Inc., Mountain View, CA
gcorrado@google.com

**Jeffrey Dean**
Google Inc., Mountain View, CA
jeff@google.com

**Abstract**

We propose two novel model architectures for computing continuous vector representations of words from very large data sets. The quality of these representations is measured in a word similarity task, and the results are compared to the previously best performing techniques based on different types of neural networks. We observe large improvements in accuracy at much lower computational cost, i.e. it takes less than a day to learn high quality word vectors from a 1.6 billion words data set. Furthermore, we show that these vectors provide state-of-the-art performance on our test set for measuring syntactic and semantic word similarities.
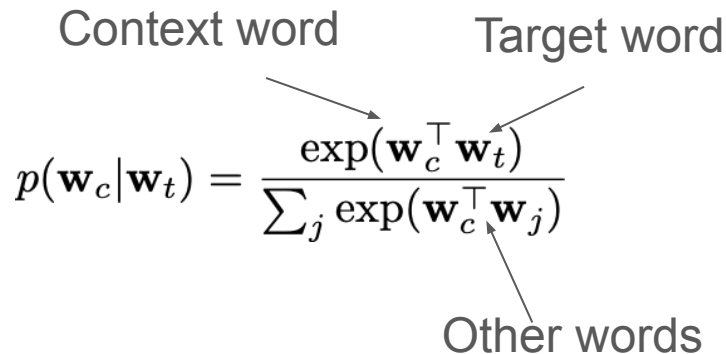
# Word2Vec - Training

**SkipGram**

$$P(w_{t-2}|w_t)$$
$$P(w_{t-1}|w_t)$$
$$P(w_{t+2}|w_t)$$
$$P(w_{t+2}|w_t)$$

| A | cup | of | coffee | is | on | the | table |
|---|-----|-----|--------|-----|-----|-----|-------|

context words in window of size 2     Center word     context words in window of size 2

**cBow**

$$P(w_t|w_{t-2})$$
$$P(w_t|w_{t-1})$$
$$P(w_t|w_{t+1})$$
$$P(w_t|w_{t+2})$$

| A | cup | of | coffee | is | on | the | table |
|---|-----|-----|--------|-----|-----|-----|-------|

context words in window of size 2     Center word     context words in window of size 2

# Looking closer…

Context word    Target word

$$p(\mathbf{w}_c|\mathbf{w}_t) = \frac{\exp(\mathbf{w}_c^\top \mathbf{w}_t)}{\sum_j \exp(\mathbf{w}_c^\top \mathbf{w}_j)}$$

(One option is to have two embeddings for each word, one as target and one as context.)

Other words

Cross-Entropy Loss:

$$\mathcal{L}_{\mathbf{W}} = - \sum_{(c,t) \in D} \log\left(p(\mathbf{w}_c|\mathbf{w}_t)\right)$$

# Word2Vec Architecture - SkipGram

Predict context word from target word!

coffee → cup



**input**

**hidden**

**output**
softmax

$0$

$1$

$0$

$d$

IVI   Matrix W

$h_1$

$h_i$

$h_d$

IVI

$d$   Matrix W'

$y_1$

$y_i$

$y_{|V|}$

Cross-Entropy
Loss

$0$

$1$

$0$

One hot
encoding of the
context word
(dim=|V|)

d-dimensional vector

IVI-dimensional vector

# Word2Vec Architecture - CBOW (continuous bag of words)

**input**

(cup, of, is, on) → coffee

**hidden**

**output**
softmax

$x_1$

$0$
·
·
·
·
$1$
·
·
·
·
$0$

$x_2$

$0$
·
·
·
·
$1$
·
·
·
·
$0$

|V| Matrix W, d

$h_1$
·
·
$h_i$
·
·
·
$h_d$

d Matrix W', |V|

$y_1$
·
·
$y_i$
·
·
·
$y_{|V|}$

d-dimensional vector
(average of vectors of all input words)

$\mathrm{avg}(xW) = h$

|V| dimensional
vector

# X 2 vec

- Generate vector representations (embeddings) for various data types
- Examples:
  - Word2Vec
  - Doc2Vec
  - Node2Vec
  - Item2Vec
  - Sent2Vec
  - Gene2Vec



Visualize: https://projector.tensorflow.org/

Explore: http://epsilon-it.utu.fi/wv_demo/

# In vector space…

# Word embeddings are time-dependent (why?)

- Semantic similarity of words depends on *time.*



**A**
daft  **gay (1900s)**
flaunting      sweet
tasteful          cheerful
                pleasant
frolicsome
witty    **gay (1950s)**
         bright
gays    bisexual
      homosexual
**gay (1990s)**
lesbian

**B**
spread
                        sow
**broadcast (1850s)** seed
circulated          sows
                scatter
      **broadcast (1900s)**
      newspapers
      television
      radio
bbc **broadcast (1990s)**

**C**
      solemn
      **awful (1850s)**
      majestic
awe
dread          pensive
            gloomy
      horrible
appalling terrible
**awful (1900s)**
            wonderful
      **awful (1990s)**
      awfully weird

# Problems with word2vec

- Words with multiple meanings only have one representation
  - eg. **bank** of river or **bank** of money
  - Need contextual information
- Limited Context
  - only trained on words within the context window

# Language Modelling History



Legend:
- Before Transformer
- After Transformer

Word2Vec N-grams — 2013

RNN / LSTM — 2014

Transformer — 2016

BERT — 2018

T5 — 2019

GPT-3 — 2020

PaLM — 2022

Bard / GPT-4 — 2023

2008

# Self-Attention



A bat flew out of the dugout, startling the baseball player and making him drop his bat.

# Self-Attention

Each word pays attention to the words around it, focusing more on the words which provide important information about its meaning.

A bat flew out of the dugout, startling the baseball player and making him drop his bat.

# Transformer Architecture

Encoder:
Computes
**contextual**
**word**
**embeddings**
from context.

Decoder:
Predicts the
next word
given
context.

Kaffee

Linear Classifier

Encoder

Decoder

xN

Feed Forward

Self-Attention

xN

Feed Forward

Cross-Attention

Masked Self-Attention

Positional
Encoding

Positional
Encoding

Embedding

Embedding

I like black coffee

<START> Ich mag schwarzen

# Self-Attention

# Self-Attention

# Discuss:

- Q, K, V are all (n x d) matrices.

- What is the shape of QK$^T$?
  - What does this matrix represent?
- What is the shape of the final output?
  - What does this matrix represent?

$$Attention(Q, K, V) = Softmax(\frac{QK^T}{d_k})V$$

# Multi-Head Attention

What if I want to pay attention to different things at the same time!?

Content-based          This is my big red dog, Clifford.

Description-based      This is my big red dog, Clifford.

Reference-based        This is my big red dog, Clifford.

What's useful depends on the task. How do I pick what to do?

# Multi-Head Attention

- The Scaled Dot-Product Attention attends to one or few entries in the input key-value pairs.
- Idea: apply Scaled Dot-Product Attention multiple times on the linearly transformed inputs.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\textbf{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

# Transformer Architecture

# Cross-Attention

# Cross-Attention



$$\mathbf{s_1} = \frac{\mathbf{q_1}\mathbf{K}^\top}{\sqrt{d}}$$

$$\boldsymbol{\alpha_1} = \mathrm{softmax}(\mathbf{s_1})$$

$$\mathbf{x_1} = \sum_{i=1}^{4} \alpha_{1,i}\mathbf{v}_i$$

softmax

Keys (K)

Queries (Q)

Values (V)

$\mathbf{W_q}$

$\mathbf{W_k}$

$\mathbf{W_v}$

Decoder Input

Encoder Output

Output Sequence

Ich mag schwarzen Kaffee

Linear Classifier

xN

Feed Forward

Self-Attention

Positional Encoding

Embedding

I like black coffee

xN

Feed Forward

Cross-Attention

Masked Self-Attention

Positional Encoding

Embedding

<START> Ich mag schwarzen

# Transformer Architecture

# Self-Attention vs. Masked Self-Attention

In masked self-attention each word only pays attention to words from the past, but not the future.

**Quiz**: Why is this advantageous in the decoder?



Self-Attention

Masked Self-Attention

# Transformer Architecture

# Point-wise Feed-forward Networks

- Purpose

  - Applies non-linear transformations to the output of the attention layer

- Equation

  - $FFN(x) = max(0, xW_1 + b_1)W_2 + b_2$

    - where W and b are learned weights and biases

- These FFN is applied separately to each position

- Followed by Layer Normalization



Layer Normalization

Sentence Length

Feature

Batch Dimension

# Transformer Architecture

# Input Embeddings

- Replace tokens with continuous vectors

- Made up of two components:

  - token embeddings

  - positional encoding depends on position and dimension index as follows (next slide)

- Word order is important:

  - "Sally stole money from John"

  - "John stole money from Sally"

# Positional Embedding

Encodes the position of each token in the sentence.





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

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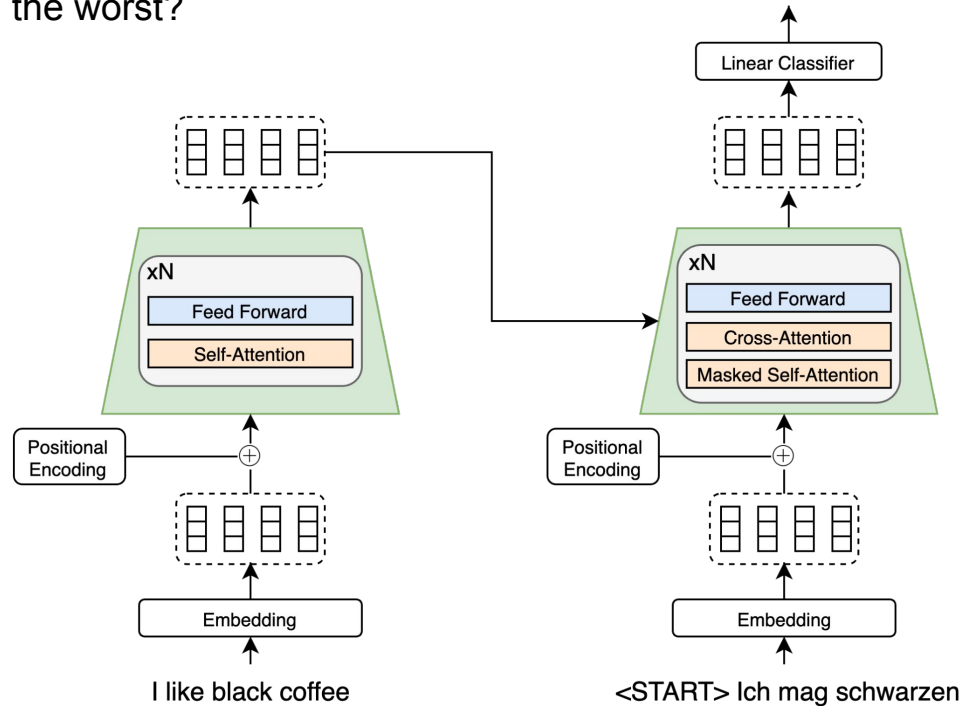
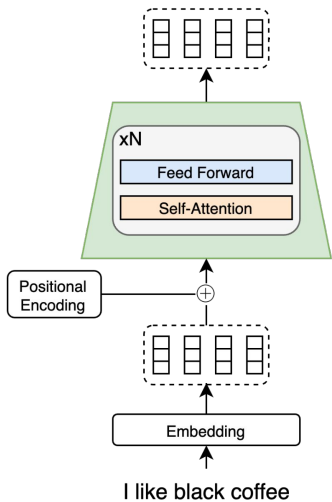# Transformer Summary

# Discuss:

- How does the transformer scale with sequence length?
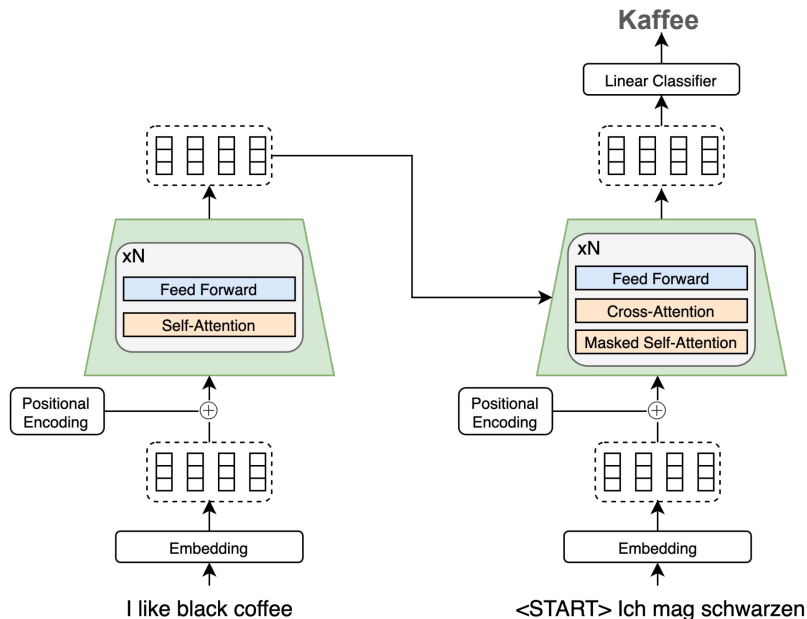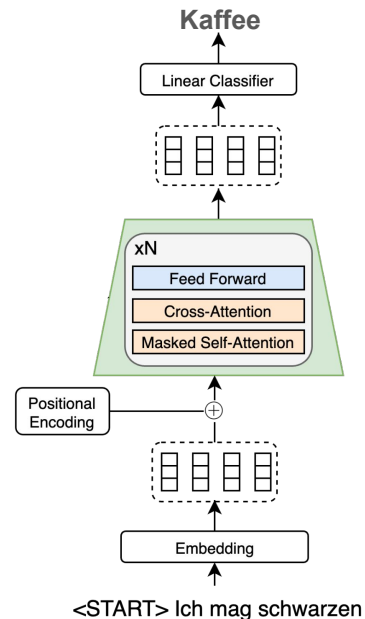  - Which part of the transformer scales the worst?

# BERT

Encoder only

# T5/BART

Encoder - Decoder

# GPT

Decoder only

**BERT**



Feed Forward

Self-Attention

xN

Positional Encoding ⊕

Embedding

I like black coffee

**T5/BART**

Feed Forward

Self-Attention

xN

Positional Encoding ⊕

Embedding

I like black coffee

Kaffee

Linear Classifier

Feed Forward

Cross-Attention

Masked Self-Attention

xN

Positional Encoding ⊕

Embedding

<START> Ich mag schwarzen

**GPT**

Kaffee

Linear Classifier

Feed Forward

Cross-Attention

Masked Self-Attention

xN

Positional Encoding ⊕
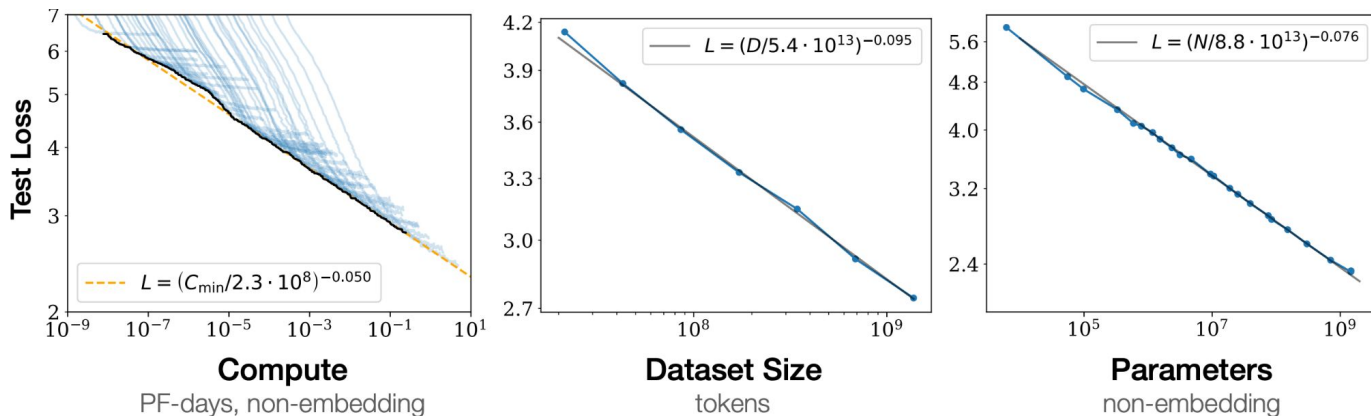
Embedding

<START> Ich mag schwarzen

# Scaling Laws

$$C = C_0 N D$$
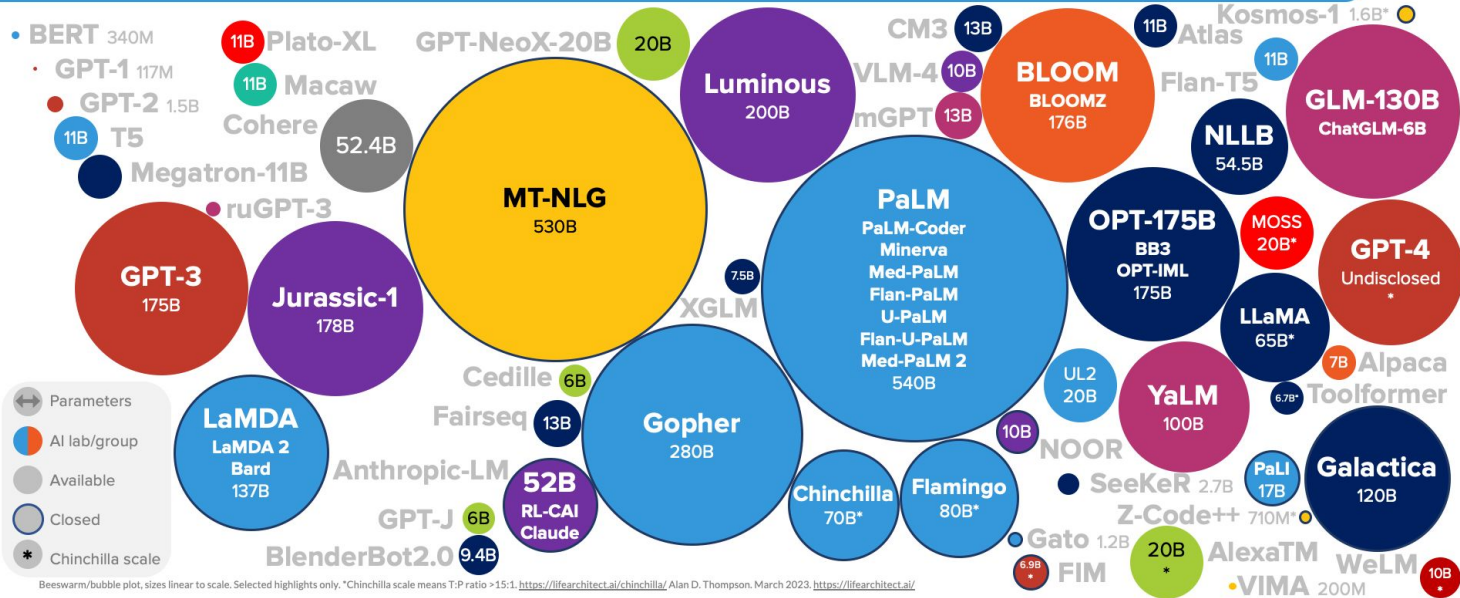$$L = \frac{A}{N^\alpha} + \frac{B}{D^\beta} + L_0$$

- Performance improves predictably with increased compute, data, and parameters
  - Can actually fit power laws!
  - Predict performance before training!



**Figure 1**  Language modeling performance improves smoothly as we increase the model size, datasetset size, and amount of compute[2] used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.
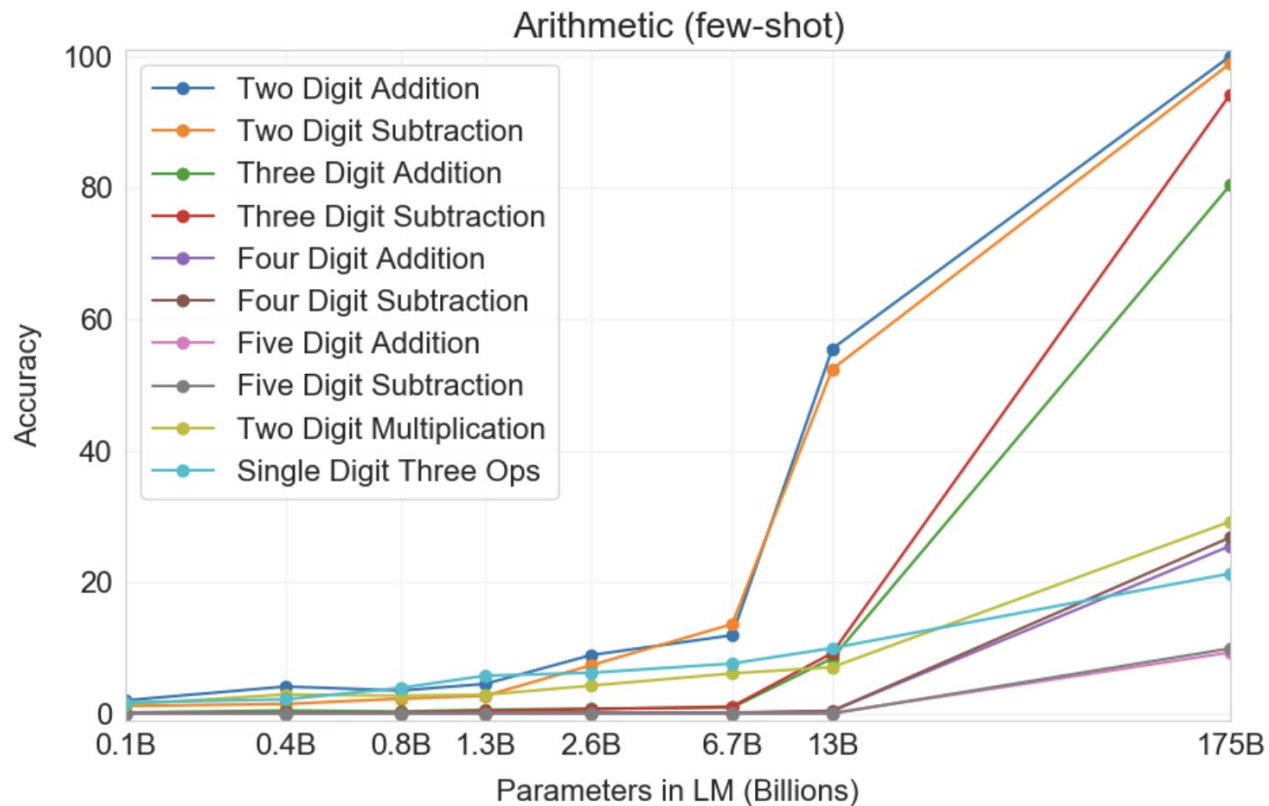
# GPT-3



**LANGUAGE MODEL SIZES TO MAR/2023**

LifeArchitect.ai/models

# Emergent Capabilities



Arithmetic (few-shot)

# Is Emergence a Mirage?

- Look at four digit addition under different metrics
  - Emergence can be an artefact of the evaluation metric
- Addition capabilities improve smoothly when using a more granular metric
  - Exact match accuracy -> Token edit distance



Schaeffer, Rylan, Brando Miranda, and Sanmi Koyejo. "Are emergent abilities of large language models a mirage?." Advances in Neural Information Processing Systems 36 (2024).