

Artificial Intelligence

CSC 665

tyler dae devlin

Course Overview

8.22.2023

Demo

Eliza vs. ChatGPT

ELIZA

- Developed in the mid 1960s
- Powered by pattern matching and substitution rules
- A beginning programmer can implement a toy version from scratch

ChatGPT

- Released 11/2022
- Powered by deep learning
- It is impossible for any individual to recreate this model from scratch
- But we will try to understand its core mechanisms by the end of the course

Introductions

Who am I?

- Tyler
- Machine learning engineering at Yelp and LinkedIn
- Applied mathematics and biology in college
- Been playing bass for 2 years



Who are you?

What this class is about

MIND

A QUARTERLY REVIEW
OF
PSYCHOLOGY AND PHILOSOPHYI.—COMPUTING MACHINERY AND
INTELLIGENCE

BY A. M. TURING

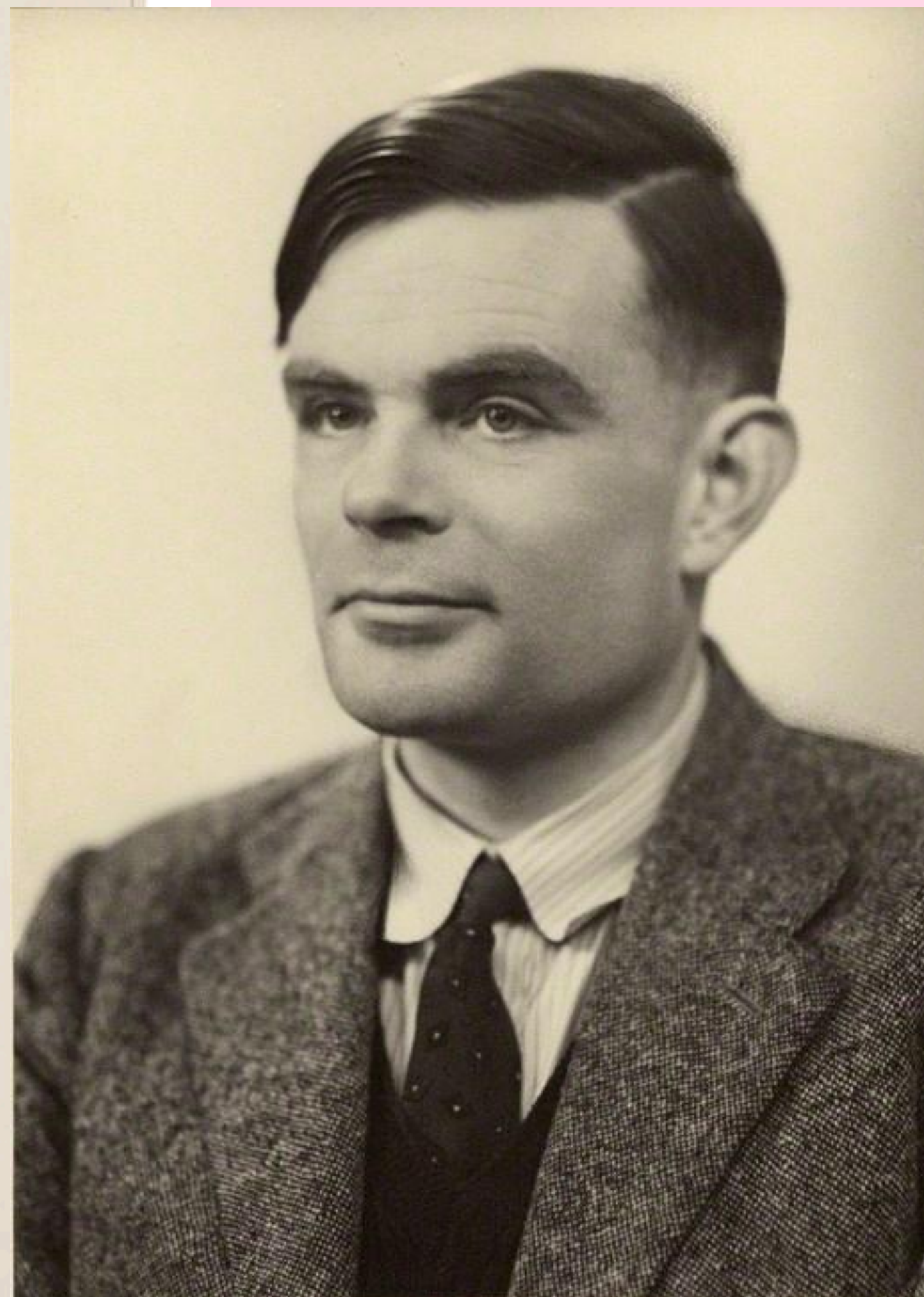
1. *The Imitation Game.*

I PROPOSE to consider the question, 'Can machines think?' This should begin with definitions of the meaning of the terms 'machine' and 'think'. The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous. If the meaning of the words 'machine' and 'think' are to be found by examining how they are commonly used it is difficult to escape the conclusion that the meaning and the answer to the question, 'Can machines think?' is to be sought in a statistical survey such as a Gallup poll. But this is absurd. Instead of attempting such a definition I shall replace the question by another, which is closely related to it and is expressed in relatively unambiguous words.

The new form of the problem can be described in terms of a game which we call the 'imitation game'. It is played with three people, a man (A), a woman (B), and an interrogator (C) who may be of either sex. The interrogator stays in a room apart from the other two. The object of the game for the interrogator is to determine which of the other two is the man and which is the woman. He knows them by labels X and Y, and at the end of the game he says either 'X is A and Y is B' or 'X is B and Y is A'. The interrogator is allowed to put questions to A and B thus:

C: Will X please tell me the length of his or her hair?
Now suppose X is actually A, then A must answer. It is A's

Where it all began



Alan Turing (1912–1954)

The Turing test

Turing machines

Turing completeness

*People have been thinking of
AI for as long as there have
been computers.*

What is special about AI?

AI versus

- programming languages
- operating systems
- databases
- networking

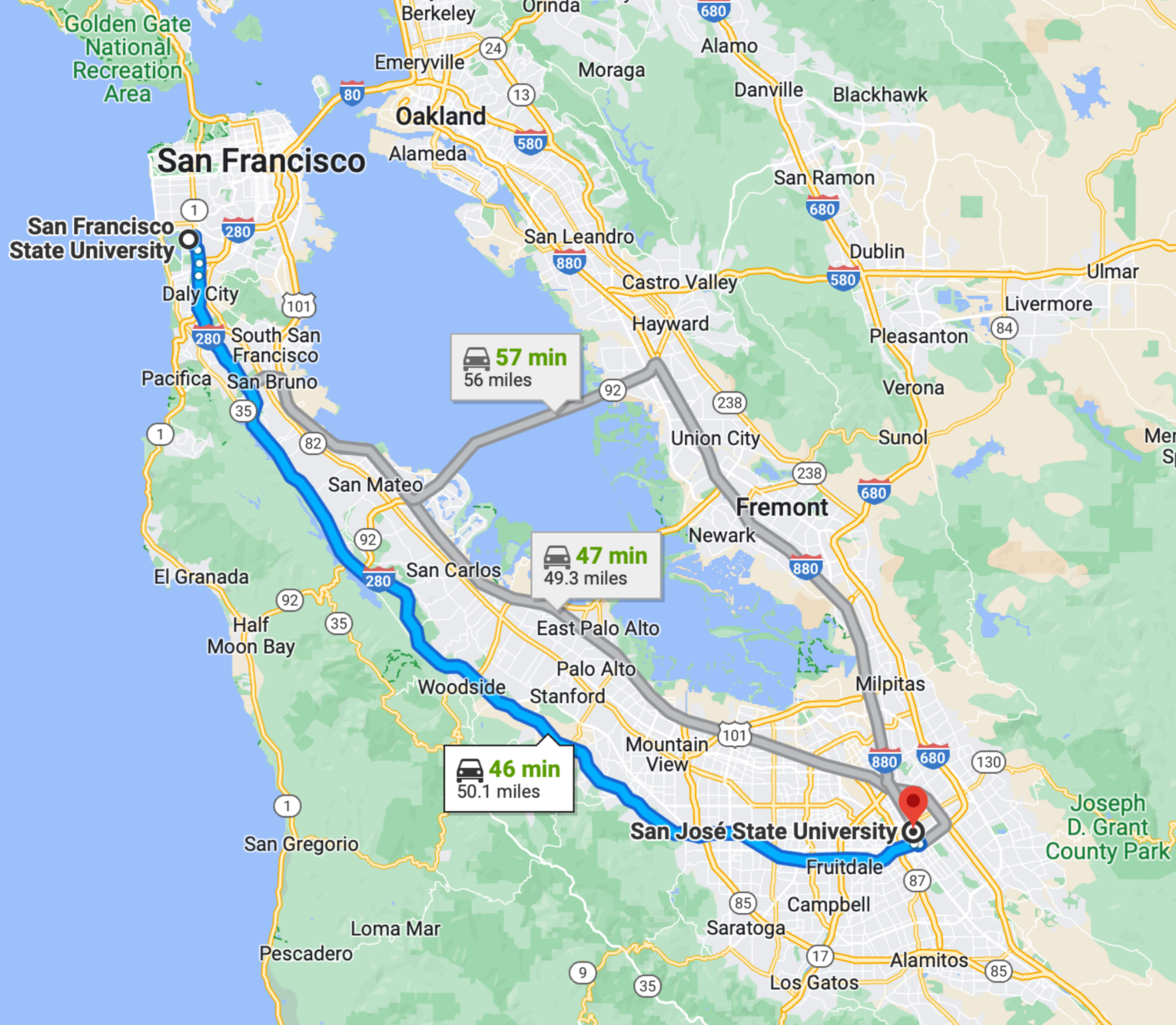
This course: a survey of some things in CS that were called AI at one point or another

What is AI?

- **Two dimensions**
 - thought vs. behavior
 - human vs. rational
- AI research has mainly focused on producing agents capable of **rational behavior**
- Agents should “do the right thing”
- Agents should be capable of behavior beyond their explicit programming
- **Example:** Is a Roomba AI?

Framework

- **Modeling and representation:** How do we build an abstract representation to model the real world? (Always a lossy process.)
- **Inference and prediction:** How do we answer a specific question given a model?
- **Learning:** How can we use data to produce better models and make better inferences?
- **Politics:** How should we allocate harms and benefits in society, and who decides?



Example: Driving Directions

(Is this AI?)

Framework

- **Model:** Use a graph to model the world — intersections are nodes, roads are edges.
- **Inference:** Use graph search algorithms to find the shortest path from start to destination.
- **Learning:** Use data to annotate edges with driving times.
- **Politics:** Road congestion, road access, car dominance, landmark prominence.
Mapmaking is always political.

Entire course on one slide

- **Search:** make decisions by looking ahead
- **Logic:** deduce new facts from existing facts
- **Constraints:** find a way to satisfy a given specification
- **Probability:** reason quantitatively about uncertainty
- **Learning:** make future predictions from past observations

Course mechanics

tddevlin.com/csc665-fall23/

Prerequisites

- Officially, CSC 413 (software development)
- Transitively, CSC 220 (data structures)
- Ideally, CSC 230 (discrete math), CSC 510 (algorithms), MATH 324 (probability)
- If you don't have the official prereq, either:
 - submit evidence of background knowledge to me via email
 - talk to me after class
- *Students who don't have the prereqs will be dropped*

Prerequisites

- **Prerequisites elsewhere**
 - [SJSU](#): DS & algorithms, OOP
 - [SCU](#): DS, discrete math
 - [Cal](#): DS, discrete math & probability
 - [Stanford](#): intro CS, discrete math, probability, linear algebra
- We will use Python for all homework assignments
- In AI, the more math you know the better
- Problems involving significant math will be optional extra credit
- Use the first homework to calibrate

Attention is All You Need,

Vaswani et al.

85K citations

The mechanism at the core of
all modern LLMs

Math content

- summation notation
- matrix operations
- big-oh notation

3.2.2 Multi-Head Attention

Instead of performing a single attention function with d_{model} -dimensional keys, values and queries, we found it beneficial to linearly project the queries, keys and values h times with different, learned linear projections to d_k , d_k and d_v dimensions, respectively. On each of these projected versions of queries, keys and values we then perform the attention function in parallel, yielding d_v -dimensional output values. These are concatenated and once again projected, resulting in the final values, as depicted in Figure 2.

Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. With a single attention head, averaging inhibits this.

⁴To illustrate why the dot products get large, assume that the components of q and k are independent random variables with mean 0 and variance 1. Then their dot product, $q \cdot k = \sum_{i=1}^{d_k} q_i k_i$, has mean 0 and variance d_k .

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

Where the projections are parameter matrices $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$ and $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$.

In this work we employ $h = 8$ parallel attention layers, or heads. For each of these we use $d_k = d_v = d_{\text{model}}/h = 64$. Due to the reduced dimension of each head, the total computational cost is similar to that of single-head attention with full dimensionality.

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. n is the sequence length, d is the representation dimension, k is the kernel size of convolutions and r the size of the neighborhood in restricted self-attention.

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

Compare with a “real” math paper

Lemma 7.4. Let $\mathrm{OG}_{(g,m,r,s,a,d)}^{\mathrm{stat},\mathrm{red}} \subset \mathrm{OG}_g^{\mathrm{stat},\mathrm{red}}$ be as in Lemma 6.6. Then

$$\int_{\mathrm{OG}_{(g,m,r,s,a,d)}^{\mathrm{stat},\mathrm{red}}} \alpha = \frac{-1}{(a_1!) \dots (a_s)!} \frac{(r+n-2)!}{r!} \cdot B_r,$$

where B_r is the Bernoulli number and $n = |a| = a_1 + \dots + a_s$.

Proof. First dispense with a special case that $(r, n) = (0, 2)$. Then either $s = 1$, $a_1 = 2$, and both sides are $-1/2$, or $s = 2$, $a_1 = a_2 = 1$, and both sides are -1 .

Otherwise, an object of $\mathrm{OG}_{(g,m,r,s,a,d)}^{\mathrm{stat},\mathrm{red}}$ has n many tails, each of which consists of an edge ending in the tail point, which is a valence-1 vertex. If we pick a total ordering of the set of the a_1 many tails with value d_1 , a total ordering of the set of the a_2 many tails with value d_2 , etc, we arrive at an object of the groupoid $\mathbb{F}_{r,n}^p$ defined exactly as $\mathbb{F}_{r,n}^\circ$, except that we do not require the marking function $m : \{1, \dots, n\} \rightarrow V$ to be injective. The superscript p stands for “pure,” as in [CGP22]. Accounting for the choices of orderings of marked points, and the n many extra edges where the tails are attached, we therefore get

$$\int_{\mathrm{OG}_{(g,m,r,s,a,d)}^{\mathrm{stat},\mathrm{red}}} \alpha = \frac{(-1)^n}{(a_1!) \dots (a_s)!} \int_{\mathbb{F}_{r,n}^p} \alpha = \frac{(-1)^n}{(a_1!) \dots (a_s)!} \sum_{(G,m) \in \mathbb{F}_{r,|a|}^p} \frac{(-1)^{|E(G)|}}{|\mathrm{Aut}(G)|}.$$

The lemma now follows from the formula

$$(18) \quad \sum_{(G,m) \in \mathbb{F}_{r,n}^p} \frac{(-1)^{|E(G)|}}{|\mathrm{Aut}(G)|} = (-1)^{n+1} \frac{(r+n-2)!}{r!} \cdot B_r,$$

valid when $2r - 2 + n > 0$. The case $n = 0$, $r \geq 2$ of this formula was given by Kontsevich [Kon93]; see [Ger04, §7.1] for a proof. The general case is easily deduced by induction on n ; see Appendix 10. \square

8. CONCLUSION

In light of Lemma 7.1, we may apply (8) to the functor $\mathcal{R} : \mathrm{OG}_g^{\mathrm{stat}} \rightarrow \mathrm{OG}_g^{\mathrm{stat},\mathrm{red}}$ to rewrite the formula in Corollary 5.9 as

$$\begin{aligned} z_g &= \int_{\mathrm{OG}_g^{\mathrm{stat}}} \alpha \cdot \beta \cdot \gamma \\ &= \int_{\mathrm{OG}_g^{\mathrm{stat},\mathrm{red}}} (\mathcal{R}_* \alpha) \cdot \beta \cdot \gamma. \end{aligned}$$

We may now replace $\mathrm{OG}_g^{\mathrm{stat},\mathrm{red}}$ by the coproduct in Lemma 6.6, and observe by Proposition 7.2 that α and β are constant functions on each $\pi_0(\mathrm{OG}_{(g,m,r,s,a,d)}^{\mathrm{stat},\mathrm{red}})$.

$$z_g = \sum_{(k,m,r,s,a,d)} \left(\prod_{i=1}^s (-\mu(m/d_i))^{a_i} \int_{\mathrm{OG}_{(g,m,r,s,a,d)}^{\mathrm{stat},\mathrm{red}}} \alpha \right) \left(P_m^{1-r} \prod_{i=1}^s \left(\frac{P_{d_i}}{P_m} \right)^{a_i} \right) \left(m^{r-1} \prod_{p|D} (1 - p^{-r}) \right).$$

Combining with the formula in 7.4 then finishes the proof of Theorem 1.1. \square

Attendance

- Attendance is optional
- Participation grade: up to 2% extra credit (quality over quantity)
- Forms of participation:
 - Attend lecture and ask questions
 - Attend lecture and answer questions
 - Post on the discussion forum
- Please *do not* come to class if you are sick

Academic integrity

- You are encouraged to study together, but any work you submit must be your own
- Talk to me if you feel yourself slipping
- *You should not expect to pass the class if you violate the honor code*

Academic integrity

- **Okay**

- Discussing assignments verbally
- Taking brief notes during discussions with classmates
- Posting high-level pseudocode on the discussion board

- **Not okay**

- Copy-pasting code
- Taking pictures/screenshots of homework solutions
- Typing on someone else's laptop
- Typing on your own laptop while looking at someone else's work

Other section

- CSC 665 is also offered on Mondays and Wednesdays, 5 – 6:15 pm
- **Instructor:** Lothar Narins
- Different approach, different style, different emphasis
- Consider checking it out

Things to do

- Homework 0 is out; due next Monday at midnight.
- Make a campuswire account and join the class (see Canvas).
- If you don't have the official prereq, either:
 - submit evidence of background knowledge to me via email
 - talk to me after class

Welcome and good luck!