

10-423/623: Generative AI Lecture 1 – Recurrent Neural Networks and Language Modelling

Henry Chai & Matt Gormley

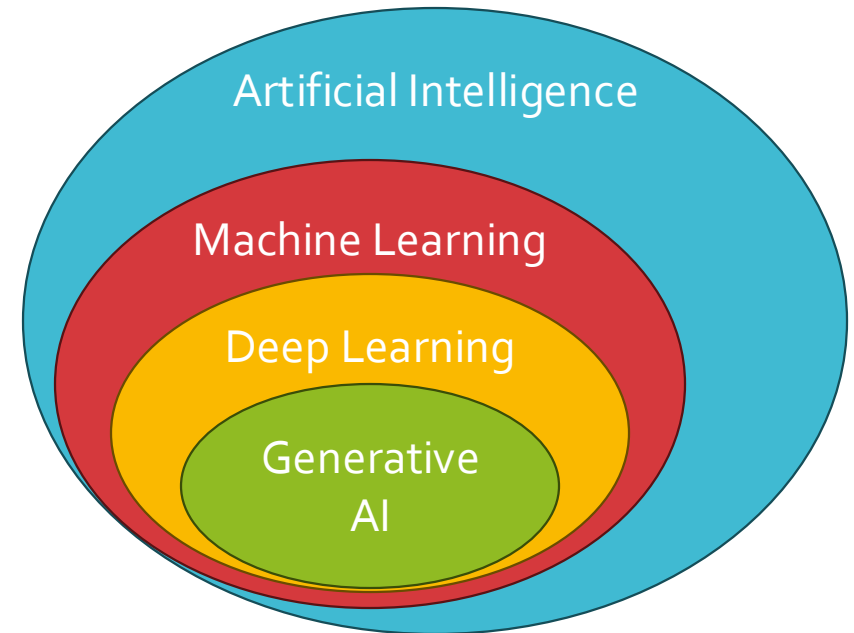
8/26/24

What is Generative AI?

What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior
- Sub-goals:
 - Perception
 - Reasoning
 - Planning
 - Communication
 - Creativity
 - Control / Motion / Manipulation
 - Learning



What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior
- Sub-goals:
 - Perception
 - Reasoning
 - Planning
 - Communication
 - Creativity
 - Control / Motion / Manipulation
 - Learning
- Large language models are capable (to varying degrees) or comprehending and responding to natural language

What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior
- Sub-goals:
 - Perception
 - Reasoning
 - Planning
 - Communication
 - Creativity
 - Text-to-image models (e.g., Midjourney) or text-to-music models (MusicGen)
 - Control / Motion / Manipulation
 - Learning

What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior

- Sub-goals:

- Perception

- Reasoning

- Planning

- Communication

- Creativity

- Control / Motion / Manipulation

- Learning

- Multimodal models can answer natural language questions about images and perform zero-shot classification

What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior
- Sub-goals:
 - Perception
 - Reasoning
 - Planning
 - Communication
 - Creativity
 - Control / Motion / Manipulation
 - Learning
- LLMs have opened new possibilities in how we approach the learning problem e.g., in-context learning

What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior
- Sub-goals:
 - Perception
 - Reasoning
 - LLMs are decent at certain reasoning tasks using e.g., chain-of-thought prompting
 - Planning
 - Communication
 - Creativity
 - Control / Motion / Manipulation
 - Learning

What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior
- Sub-goals:
 - Perception
 - Reasoning
 - Planning
 - Communication
 - Creativity
 - Control / Motion / Manipulation
 - Learning
- LLM-Planner: a modified LLM that can help embodied agents (i.e., robots) perform complex, multi-step tasks by generating high-level plans

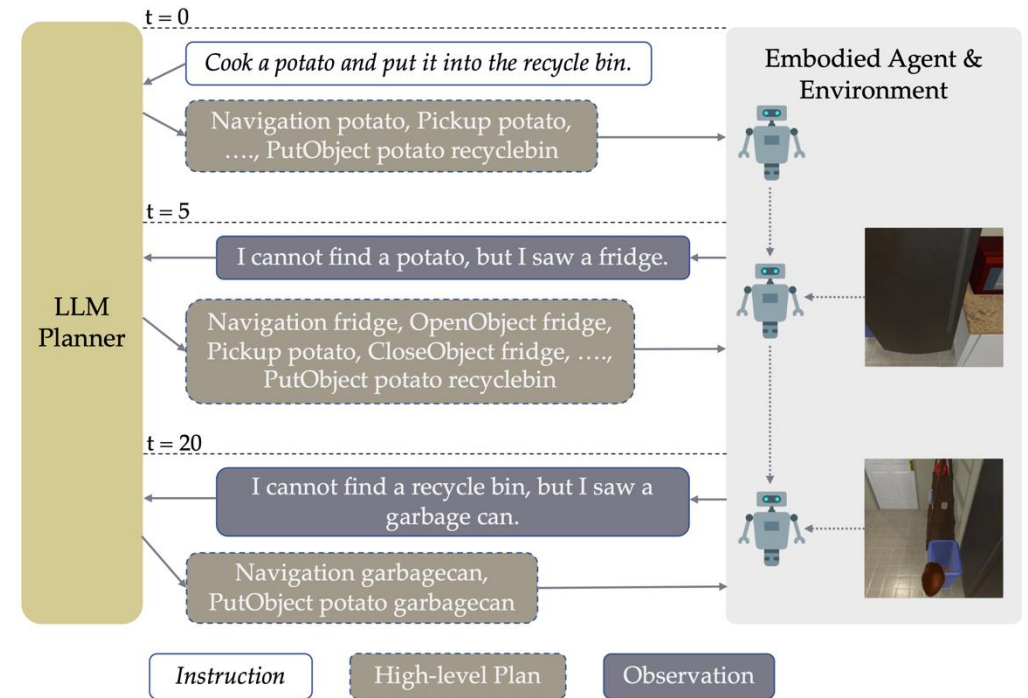
What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior

- Sub-goals:

- Perception
- Reasoning
- Planning
- Communication
- Creativity
- Control / Motion / Manipulation
- Learning



What is

AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior
- Sub-goals:
 - Perception
 - Reasoning
 - Planning
 - Communication
 - Creativity
 - Control / Motion / Manipulation
 - Learning
- DayDreamer: learns a "World Model" for reinforcement learning without the need for simulated data; allows for rapid training of RL agents

What is AI?

- The field of artificial intelligence is concerned with developing computer systems that mimic (or exceed) human behavior
- Sub-goals:
 - Perception
 - Reasoning
 - Planning
 - Communication
 - Creativity
 - Control / Motion / Manipulation
 - Learning

The ability to create or generate novel ideas/concepts (to “imagine”) is central to how humans successfully perform these tasks!

What is Generative AI capable of?

Teach me about narwhals in the style of a Shakespearean sonnet



Behold the narwhal, monarch of the sea,
With spiral tusk that doth the waves adorn,
A unicorn of ocean's mystery,
By nature's hand so wondrously born.

Beneath the icy realms where cold winds blow,
This creature glides through waters dark and deep,
Its ivory lance, a silent, graceful show,
A secret kept where ancient legends sleep.

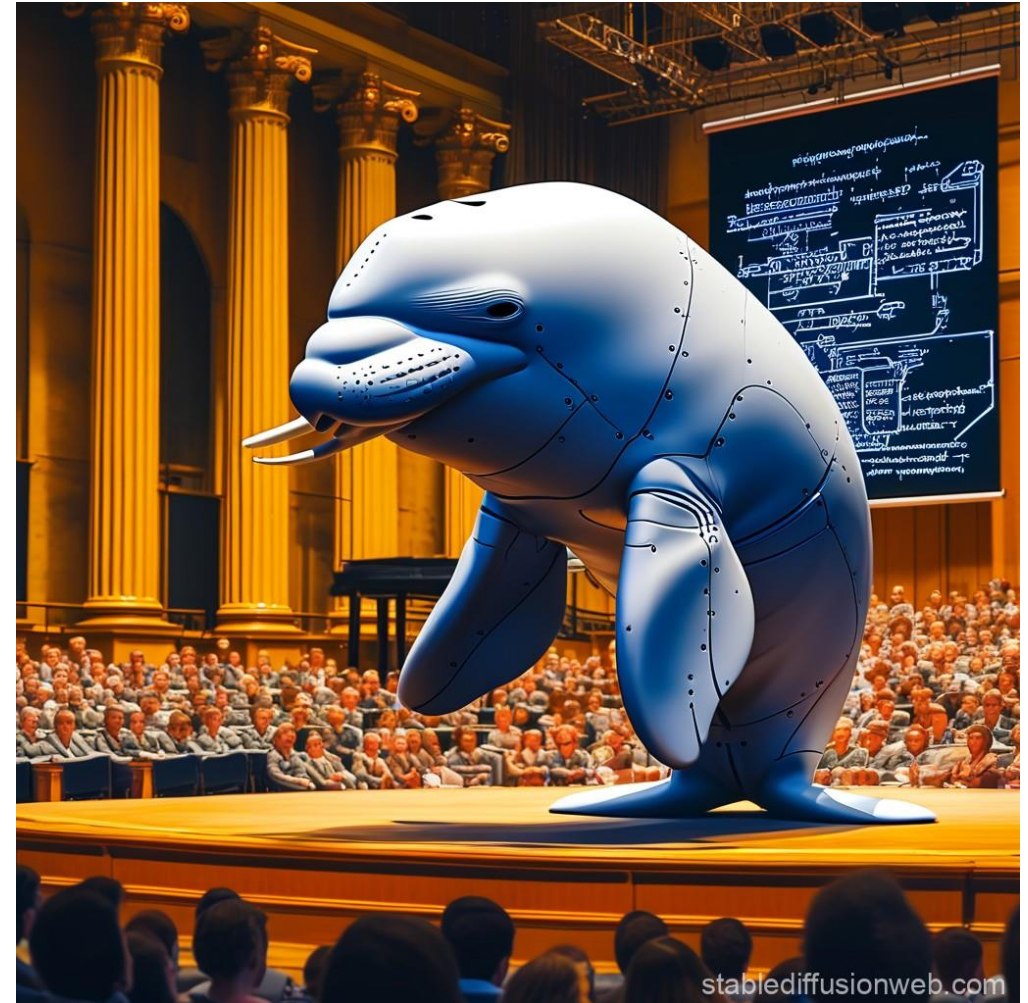
Yet not for battle is this horn designed,
But rather, like a herald in the night,
It senses all that swims or stirs behind,
Guiding its path through realms devoid of light.

O narwhal, gentle giant of the brine,
In thee, both strength and beauty intertwine.



What is Generative AI capable of?

Prompt: a picture of a narwhal giving a lecture on generative AI. Put the narwhal on a stage in front of a large audience. The narwhal should appear scholarly and is presenting off of a slide deck, gesturing at some dense equations on the slide



What is Generative AI capable of?



What is Generative AI capable of?

Whimsical Indie Folk acoustic guitar and ukulele accompaniment, playful and quirky, with a laid-back tempo around 90 BPM.



Epic cinematic score with a heavy, pulsing bassline and energetic trap beats, blending elements of electronic and orchestral music, in a 140-150 BPM tempo.



What is Generative AI capable of?

GPT-4 vs ChatGPT

Prompt:

I want to write a customized optimizer in pytorch, the optimizer should do:

1. the parameters for the optimizer is learning rate, momentum ,weigh decay, d_dim, k and factor alpha
2. the optimizer goes through each parameter in the network, and then

- Get the gradient of this parameter, reshape it to 2D by looping through the dimensions, and keeping the dimension whose actual dimension (shape) is equal to d_dim, do an SVD decomposition of the parameter so $W = U\Sigma V$.
- Keep the top k singular vectors of W, so we get a new matrix $W_k = U \Sigma_k V$.
- Keep the top 2k singular vectors of W, so we get a new matrix $W_{2k} = U \Sigma_{2k} V$.
- Normalize W_k according to the momentum of the F-norm of W_{2k} .
- Truncate each coordinate of W_k to the top alpha percentile of all the coordinates, in absolute value.
- Apply momentum on W_k , and then update the network parameters using this momentum.
- Apply weight decay.

GPT-4:

```
def step(self):
    for group in self.param_groups:
        for p in group['params']:
            if p.grad is None:
                continue
            grad = p.grad.data
            # reshape grad to 2D
            grad_2d = grad.reshape(-1, group['d_dim'])
            # SVD decomposition
            U, S, V = torch.svd(grad_2d)
            # keep top k singular vectors
            G_k = U[:, :group['k']] @ torch.diag(S[:group['k']]) @
            V[:, :group['k']].T
            # keep top 2k singular vectors
            G_2k = U[:, :2*group['k']] @ torch.diag(S[:2*group['k']])
            @ V[:, :2*group['k']].T
            # normalize G_k
            G_k /= torch.norm(G_2k, p=2)
            # truncate G_k
            threshold = torch.quantile(torch.abs(G_k),
            group['alpha'])
            G_k[torch.abs(G_k) < threshold] = 0
            # apply momentum
            if 'momentum_buffer' not in group:
                group['momentum_buffer'] = torch.zeros_like(G_k)
            group['momentum_buffer'] = group['momentum'] *
            group['momentum_buffer'] + G_k
            # update parameters
            p.data.add_(-group['lr'] * group['momentum_buffer'])
            # apply weight decay
            if group['weight_decay'] != 0:
                p.data.add_(-group['weight_decay'] * group['lr'],
                p.data)
```

ChatGPT:

```
def step(self):
    for group in self.param_groups:
        for p in group['params']:
            if p.grad is None:
                continue
            grad = p.grad.data
            grad = grad.reshape(-1, grad.shape[-1]) # reshape to 2D
            # SVD decomposition
            u, s, v = torch.svd(grad)
            # Keep top k and 2k singular vectors
            k = group['k']
            u_k = u[:, :k]
            s_k = s[:k]
            v_k = v[:, :k]
            u_2k = u[:, :2*k]
            s_2k = s[:2*k]
            v_2k = v[:, :2*k]
            # Normalize according to momentum of F-norm of G_2k
            grad_k = u_k @ torch.diag(s_k) @ v_k.t()
            norm_2k = s_2k.norm()
            grad_k = grad_k * group['momentum'] / norm_2k
            # Truncate to top alpha percentile
            alpha = group['alpha']
            percentile = grad_k.abs().flatten().kthvalue(int(alpha *
            grad_k.numel()))
            grad_k[grad_k.abs() < percentile] = 0
            # Apply weight decay
            p.data.add_(-group['lr'], grad_k + group['weight_decay']
            * p.data)
```

Figure 3.4: GPT-4 vs ChatGPT at implementing a customized Pytorch optimizer module. We highlight the **astuteness** of GPT-4 vs the **mistake** by ChatGPT at 'applying momentum'.

What is Generative AI capable of?

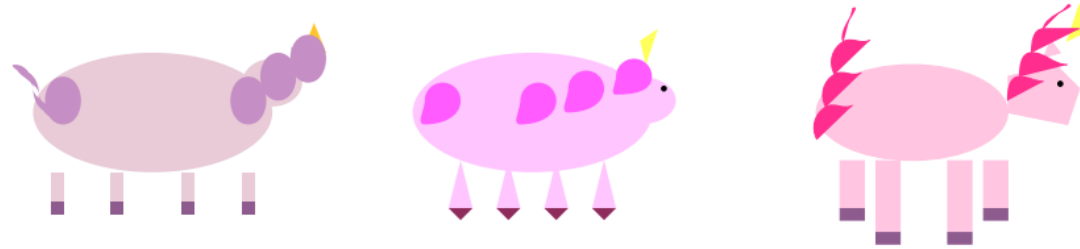


Figure 1.3: We queried GPT-4 three times, at roughly equal time intervals over the span of a month while the system was being refined, with the prompt “Draw a unicorn in TikZ”. We can see a clear evolution in the sophistication of GPT-4’s drawings.

GPT-4 vs ChatGPT

Prompt:

A function f is called *increasing* if $f(s) \leq f(t)$ holds for every $s \leq t$. Consider the formula:
 $\lim_{x \rightarrow 0} \lim_{n \rightarrow \infty} \left[\sum_{i=1}^n (f(nx) - \int_0^x \sin(nx) dx), \frac{(x + \frac{2}{8})^2 - 5}{4x - 1/x^2}, \dots, \text{partial } f(x/n) \right]$
 $\text{over partial } x * (e^{-i * \omega n x} + 3(x \# 7))$, f inverse of $(x \bmod n)$]
 multiply $\text{matrix}\{1, 1, \dots, 1 \ \backslash \ 1, 1, \dots, 1\}$
 $\{\text{grey_color}\{\text{this gives us a good impression about the function}\}\}$

Convert the above paragraph to LaTeX.

GPT-4:

A function f is called *increasing* if $f(s) \leq f(t)$ holds for every $s \leq t$. Consider the formula:

$$\lim_{x \rightarrow 0} \lim_{n \rightarrow \infty} \left[\sum_{i=1}^n (f(nx) - \int_0^x \sin(nx) dx), \frac{(x + \frac{2}{8})^2 - 5}{4x - 1/x^2}, \dots, \text{partial } f(x/n) \right] \times \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

this gives us a good impression about the function

Why is
Generative
AI worth
understanding
deeply?



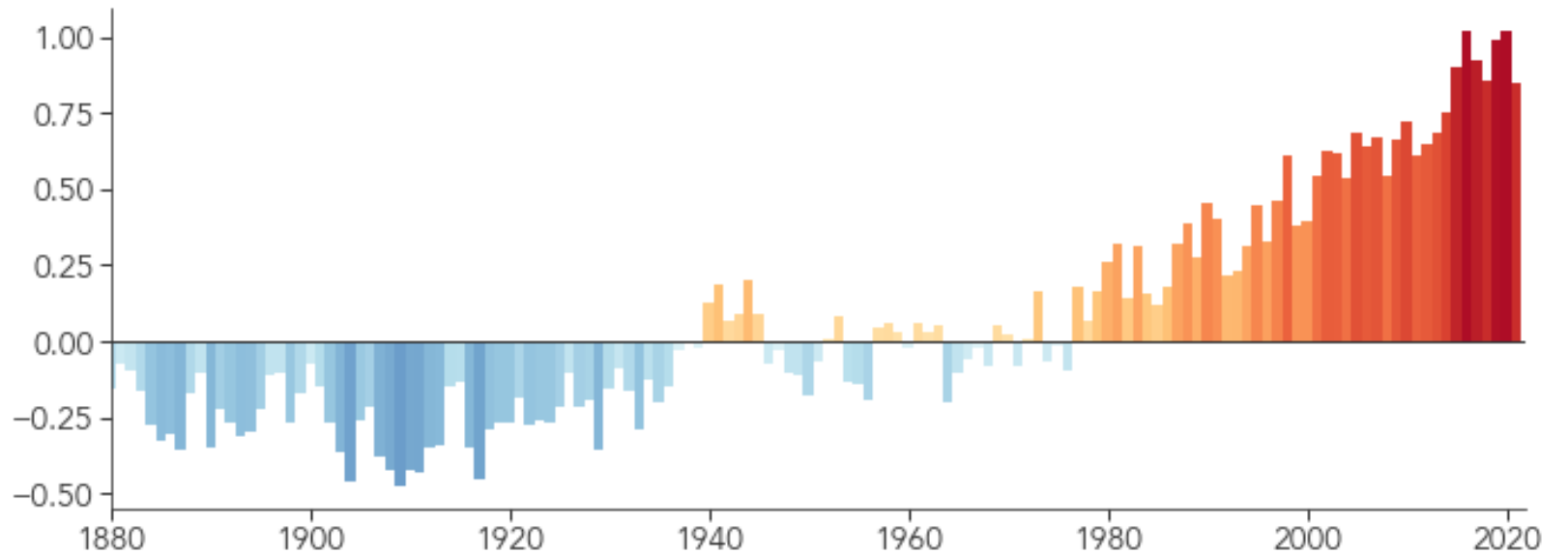
Why is Generative AI worth understanding deeply?



Why is Generative AI worth understanding deeply?

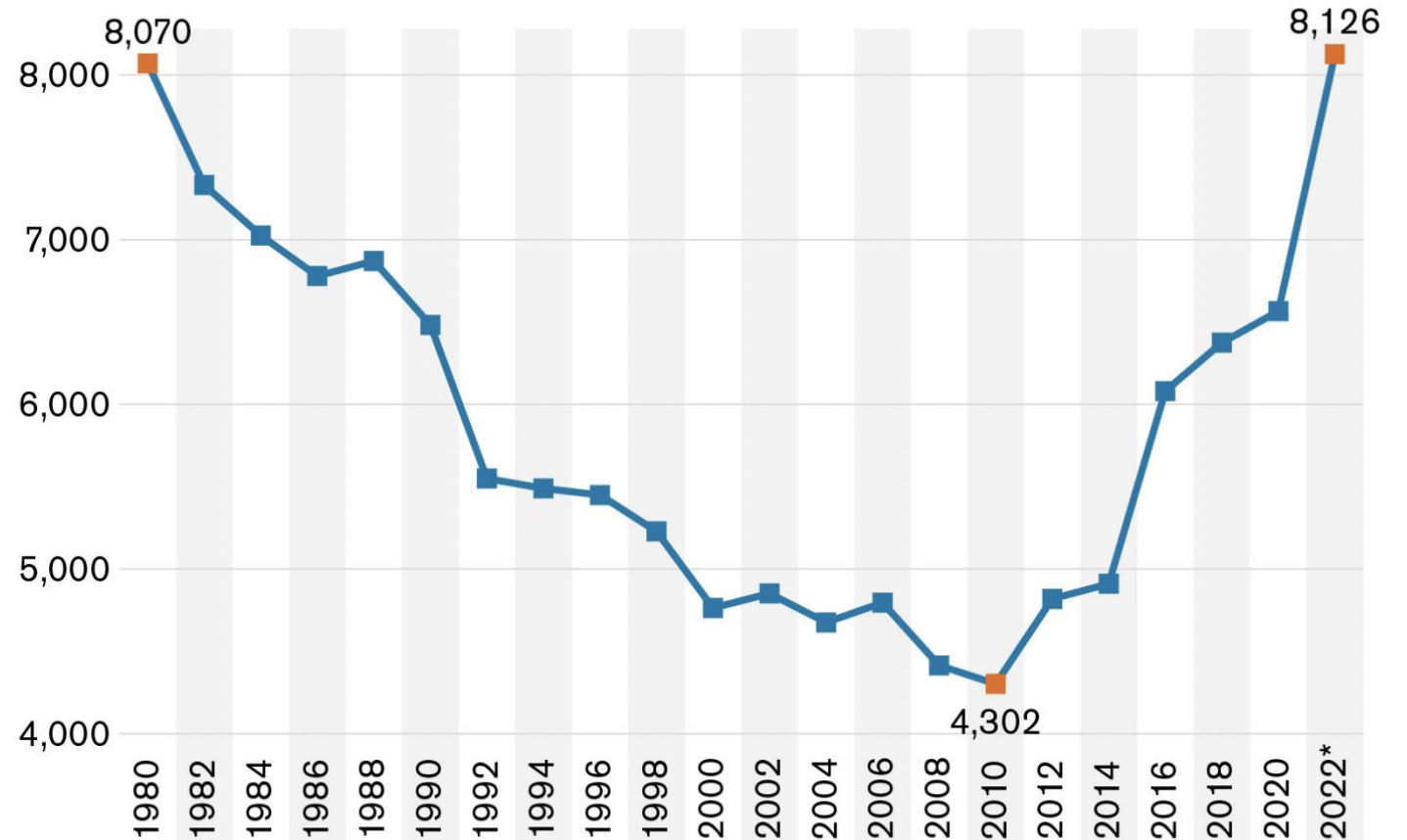
2021 ties 2018 for Sixth Warmest Year on Record

Global Temperature Anomaly (°C compared to the 1951-1980 average)



Why is Generative AI worth understanding deeply?

Figure 5 Number of Annual U.S. Pedestrian Fatalities, 1980-2022



Why is Generative AI worth understanding deeply?



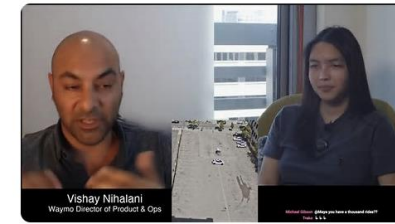
Why is Generative AI worth understanding deeply?

8/26/24

engadget · 5d · on MSN

Waymo director says the company's cars won't honk at each other anymore

Waymo's self-driving cars no longer honk when near each other, Waymo's Director of Product and Operations Vishay Nihalani ...



ars TECHNICA · 12d

Self-driving Waymo cars keep SF residents honking at each other

Silicon Valley's latest disruption? Your sleep schedule. On reported that San Francisco's South of ...

ars PC  [+11 relevant news](#) >

NEWS · 13d · on MSN

Waymo cars honk at each other throughout SF neighborhoods

San Francisco resident Randol White says he heard the noise weeks ago -- he was woken up around ...

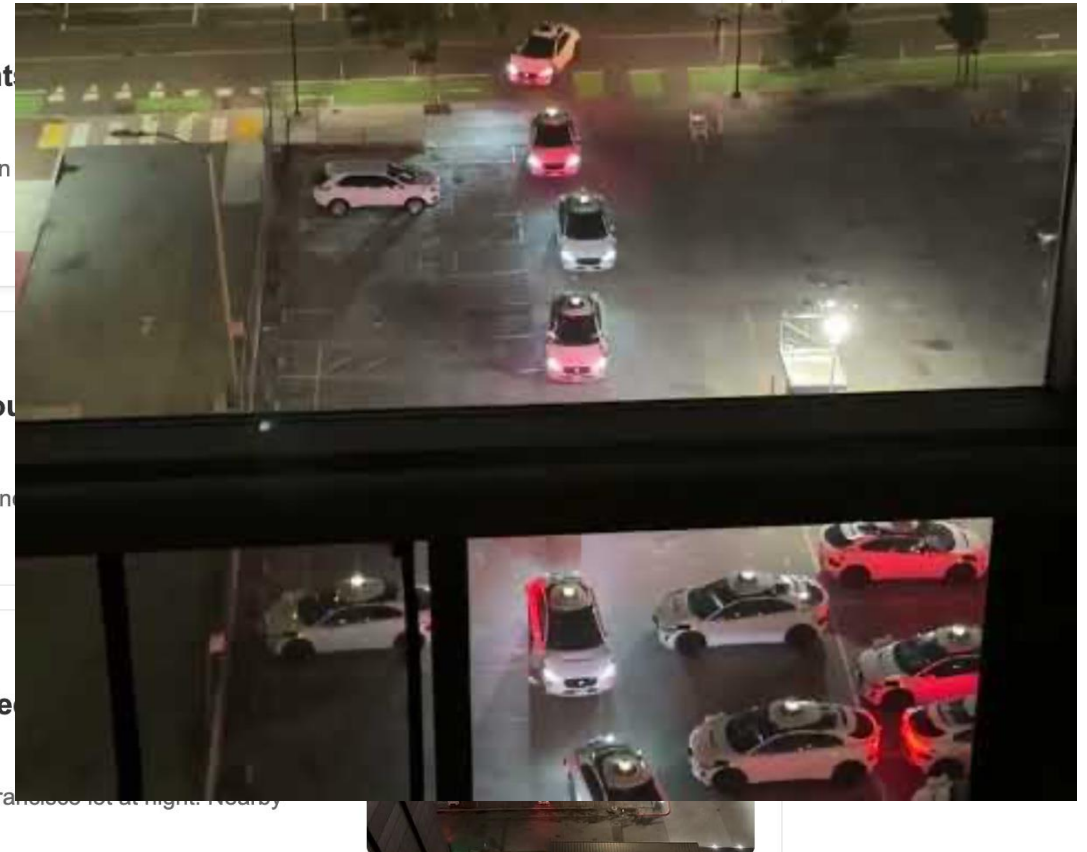
Los Angeles Times · 13d · on MSN

Driverless Waymo taxis disturb peace, sleep honking — at each other

Waymo's robotaxis have started occupying a new San Francisco lot at night, nearby residents are finding the self-driving cars ...

yahoo! · 5d

Driverless Waymo cars still honking despite software fix



Source: <https://www.bing.com/news/search?q=waymo+cars+honking+at+each+other&FORM=HDRSC7>

What is Generative AI 10-423/623: Topics

- Generative models of text
 - RNN LMs / Autodiff
 - Transformer LMs
 - Pre-training, fine-tuning, evaluation, decoding
- Generative models of images
 - CNNs / Transformers for vision
 - GANs, Conditional GANs
 - VAEs and Diffusion models
- Applying and adapting foundation models
 - Reinforcement learning with human feedback (RLHF)
 - Parameter-efficient fine tuning
 - In-context learning for text
 - In-context learning for vision
- Multimodal foundation models
 - Text-to-image generation
 - Aligning multimodal representations
 - Visual-language foundation models
- Scaling models
 - Efficient decoding strategies
 - Distributed training
 - Scaling laws and data
- What can go wrong?
 - Safety/bias/fairness
 - Hallucinations
 - Adversarial attacks
 - Cheating – how to watermark, Legal issues, e.g., copyright,...
 - Drift in performance
 - Data contamination
- Advanced Topics
 - State space models
 - Code generation
 - Audio understanding and synthesis
 - Video synthesis

What is ~~Generative~~ AI 10-423/623: Website

- <https://www.cs.cmu.edu/~mgormley/courses/10423/>
- Also <http://423.mlcourse.org> or <http://623.mlcourse.org>

What is ~~Generative~~ AI 10-423/623: Syllabus

- <https://www.cs.cmu.edu/~mgormley/courses/10423/syllabus.html>
- This entire page is **required** reading

What is ~~Generative~~ AI 10-423/623: Syllabus

- <https://www.cs.cmu.edu/~mgormley/courses/10423/syllabus.html>
- Assessments:
 - 40% homework
 - 10% quizzes
 - 20% exam
 - 25% project
 - 5% participation

What is ~~Generative~~ AI 10-423/623: Homework

- <https://www.cs.cmu.edu/~mgormley/courses/10423/coursework.html>

HW #	Topic	Model	Application	Type
0	PyTorch Primer	image classifier + text classifier	Vision + language	written + programming
1	LLMs	Transformer with RoPE and GQA	Character-level text generation	written + programming
2	Image Generation	Diffusion model	Infilling	written + programming
3	Adapters for LLMs	Llama with LoRA	Code	written + programming
4	Multimodal Foundation Models	Text-to-image model	Vision <i>and</i> language	written + programming
623 (623 students only)	Recent Research	N/A	Generative AI	video presentation

What is Generative AI 10-423/623: Homework

- <https://www.cs.cmu.edu/~mgormley/courses/10423/coursework.html>

[Home](#)

[FAQ](#)

[Syllabus](#)

[People](#)

[Schedule](#)

[Office Hours](#)

[Coursework](#)

[Links](#) ▾

Generative AI

10-423 + 10-623, Fall 2024
School of Computer Science
Carnegie Mellon University

Assignments

There will be 5 homework assignments (and a special extra assignment for 10-623 only). The assignments will consist of both theoretical and programming problems. Homework assignments will be released via a Piazza announcement explaining where to find the handout, LaTeX template, etc.

- Homework 0: PyTorch Primer
- Homework 1: Large Language Models
- Homework 2: Image Generation
- Homework 3: Applying and Adapting LLMs
- Homework 4: Multimodal Foundation Models
- Homework 623: (10-623 only)

What is ~~Generative~~ AI 10-423/623: Quizzes

- <https://www.cs.cmu.edu/~mgormley/courses/10423/coursework.html>

Home

FAQ

Syllabus

People

Schedule

Office Hours

Coursework

Links ▾

Quizzes

There will be 5 quizzes.

- Quiz 1 (Lectures 1 - 3)
- Quiz 2 (Lectures 4 - 8)
- Quiz 3 (Lectures 9 - 11)
- Quiz 4 (Lectures 12 - 15)
- Quiz 5 (Lectures 16 - 20)

- Quizzes will be held in-class, before the start of lecture

What is Generative AI 10-423/623: Quizzes

8/26/24

- <https://www.cs.cmu.edu/~mgormley/courses/10423/schedule.html>

Home

FAQ

Syllabus

People

Schedule

Office Hours

Coursework

Links ▾

Important Notes

This schedule is **tentative** and subject to change. Please check back often.

Tentative Schedule

Date	Lecture	Readings	Announcements
Generative models of text			
Mon, 26-Aug	Lecture 1 : RNN LMs / Autodiff	<ul style="list-style-type: none">• <i>Sequence Modeling: Recurrent and Recursive Nets</i>. Ian Goodfellow, Yoshua Bengio, & Aaron Courville (2016). Deep Learning, Chapter 10.1-10.5.• <i>A Framework for the Cooperation of Learning Algorithms</i>. Léon Bottou, Patrick Gallinari (1991). Advances in Neural Information Processing Systems.• <i>PyTorch: An Imperative Style, High-Performance Deep Learning Library</i>. Paszke et al. (2019). Advances in Neural Information Processing Systems.	
Wed, 28-Aug	Lecture 2 : Transformer LMs	<ul style="list-style-type: none">• <i>The Long Short-Term Memory and Other GatedRNNs</i>. Ian Goodfellow, Yoshua Bengio, & Aaron Courville (2016). Deep Learning, Chapter 10.10-10.12.• <i>Attention Is All You Need</i>. Vaswani et al. (2017). NeurIPS.• <i>The Illustrated Transformer</i>. Alammar (2018).	HW0 out
Fri, 30-Aug	Recitation: HW0		
Mon, 2-Sep	(Labor Day - No Class)		

What is Generative AI 10-423/623: Quizzes

- <https://www.cs.cmu.edu/~mgormley/courses/10423/schedule.html>

Home

FAQ

Syllabus

People

Schedule

Office Hours

Coursework

Links ▾

Wed, 11-Sep	Lecture 5 : Encoder-only Transformers / Vision Transformers	<ul style="list-style-type: none"> • <i>BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.</i> Devlin et al. (2018). NAACL. • <i>An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale.</i> Dosovitskiy et al. (2021). ICLR. 	Quiz 1 (in-class, L1-L4)
Fri, 13-Sep	Recitation: HW1		
Mon, 16-Sep	Lecture 6 : Generative Adversarial Networks (GANs) / PGM	<ul style="list-style-type: none"> • <i>Generative Adversarial Nets.</i> Goodfellow et al. (2014). NeurIPS. • <i>NeurIPS 2016 Tutorial: Generative Adversarial Networks.</i> Goodfellow (2017). NeurIPS Tutorials. 	
Wed, 18-Sep	Lecture 7 : Variational Autoencoders (VAEs) / Diffusion models (Part I)	<ul style="list-style-type: none"> • <i>Variational Inference: A Review for Statisticians.</i> Blei, Kucukelbir, and McAuliffe (2018). • <i>High-Level Explanation of Variational Inference.</i> Jason Eisner (2011). • <i>Tutorial on Variational Autoencoders.</i> Carl Doersch (2016). 	
Fri, 20-Sep	(No Recitation)		
Mon, 23-Sep	Lecture 8 : Diffusion models (Part II)	<ul style="list-style-type: none"> • <i>Deep Unsupervised Learning using Nonequilibrium Thermodynamics.</i> Sohl-Dickstein et al. (2015). MLR. • <i>Denoising Diffusion Probabilistic Models.</i> Ho et al. (2020). NeurIPS. 	HW1 due HW2 out (L4-L8)
Applying and adapting foundation models			
Wed, 25-Sep	Lecture 9 : In-context learning for text & for vision		Quiz 2 (in-class, L4-L8)

What is Generative AI 10-423/623: Exam

- <https://www.cs.cmu.edu/~mgormley/courses/10423/schedule.html>

Home

FAQ

Syllabus

People

Schedule

Office Hours

Coursework

Links ▾

Mon, 4-Nov	Lecture 18 : Efficient decoding strategies		HW4 due HW623 out
Advanced Topics			
Wed, 6-Nov	Lecture 19 : Long Context in LLM		
Fri, 8-Nov	(No Recitation)		
Mon, 11-Nov	In-Class Exam		Pass/no pass & withdrawal deadline
Wed, 13-Nov	Lecture 20 : State Space Models		
Fri, 15-Nov	(No Recitation)		Project proposal due

What is Generative AI 10-423/623: Project

- <https://www.cs.cmu.edu/~mgormley/courses/10423/schedule.html>

Home

FAQ

Syllabus

People

Schedule

Office Hours

Coursework

Links ▾

Mon, 4-Nov	Lecture 18 : Efficient decoding strategies		HW4 due HW623 out
Advanced Topics			
Wed, 6-Nov	Lecture 19 : Long Context in LLM		
Fri, 8-Nov	(No Recitation)		
Mon, 11-Nov	In-Class Exam		Pass/no pass & withdrawal deadline
Wed, 13-Nov	Lecture 20 : State Space Models		
Fri, 15-Nov	(No Recitation)		Project proposal due

What is Generative AI 10-423/623: Project

- <https://www.cs.cmu.edu/~mgormley/courses/10423/schedule.html>
- Apply generative AI to a problem / setting of your choice
- Report back to the class during a *poster session* to be held **sometime over finals period**
- You must complete the project in **groups of 3**



Prompt to ChatGPT-4o: *Create an image of three Scottish terriers in traditional Scottish outfits working collaboratively on a project for a generative AI course*

What is ~~Generative~~ AI 10-423/623: Syllabus

- <https://www.cs.cmu.edu/~mgormley/courses/10423/syllabus.html>
- Assessments:
 - 40% homework
 - 10% quizzes
 - 20% exam
 - 25% project
 - 5% participation
 - Participation will be based on a few (out-of-class) surveys/polls and some project deliverables

What is ~~Generative~~ AI 10-423/623: Prerequisites

- <https://www.cs.cmu.edu/~mgormley/courses/10423/syllabus.html>
- You must have taken either:
 - an introduction to machine learning course (e.g., 10-301 / 601 or 10-315 or 10-701 or 10-715) *or*
 - an introduction to deep learning course (e.g., 11-485 / 685 / 785)
- You **do not** need to have familiarity with deep learning or PyTorch

What is ~~Generative~~ AI 10-423/623: Collaboration

- <https://www.cs.cmu.edu/~mgormley/courses/10423/syllabus.html>
- Collaboration is strongly encouraged in this course!
- But collaboration must always be documented
- You must always write your own code/solutions
 - No re-use of found code or code written for previous assignments
- Please see the syllabus for
 - recommended collaboration practices and
 - penalties for violations
- Note the 10-423/623 collaboration policy is different (more lenient) than our 10-301/601 policies

What is Generative AI 10-423/623: Technologies

- <https://www.cs.cmu.edu/~mgormley/courses/10423/>

Home

FAQ

Syllabus

People

Schedule

Office Hours

Coursework

Links ▾

Piazza

Gradescope

Livestream

Video Recordings

OH Queue

- Piazza (for Q&A)
- Gradescope (for submitting homework)
- Zoom (for livestreams)
- Panopto (for recordings)
- Google forms (for surveys/polls)

What is Generative AI 10-423/623: Office Hours

- <https://www.cs.cmu.edu/~mgormley/courses/10423/officehours.html>

Home
FAQ
Syllabus
People
Schedule
Office Hours
Coursework
Links ▾

Generative AI

10-423 + 10-623, Fall 2024
School of Computer Science
Carnegie Mellon University

10-423/10-623 Office Hours

Today ◀ ▶ Aug 25 – 31, 2024 ▾

Week Month Ag

	Sun 8/25	Mon 8/26	Tue 8/27	Wed 8/28	Thu 8/29	Fri 8/30	Sat 8/31
12am							
1am							
2am							
3am							
4am							
5am							
6am							
7am							

What is Generative AI 10-423/623: Learning Objectives

You should be able to...

1. Differentiate between different mechanisms of learning such as parameter tuning and in-context learning.
2. Implement the foundational models underlying modern approaches to generative modeling, such as transformers and diffusion models.
3. Apply existing models to real-world generation problems for text, code, images, audio, and video.
4. Employ techniques for adapting foundation models to tasks such as fine-tuning, adapters, and in-context learning.
5. Enable methods for generative modeling to scale-up to large datasets of text, code, or images.
6. Use existing generative models to solve real-world discriminative problems and for other everyday use cases.
7. Analyze the theoretical properties of foundation models at scale.
8. Identify potential pitfalls of generative modeling for different modalities.
9. Describe societal impacts of large-scale generative AI systems.

What is Generative AI?

- Generative AI is just *modelling* and *sampling* from probability distributions that look like

$$p(x_{t+1}, x_t, \dots, x_1)$$

or

$$p(x_{t+1} | x_t, \dots, x_1)$$

Language Models

1. Convert raw text into a sequence of words (or “tokens”)

$$\mathbf{x} = [x_1, x_2, \dots, x_t]$$

2. Learn (or approximate) a joint probability distribution over sequences

$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_t)$$

3. Sample from the implied conditional distribution to generate new words (or “tokens”)

$$p(x_{t+1} | x_1, x_2, \dots, x_t) = \frac{p(x_1, x_2, \dots, x_t, x_{t+1})}{p(x_1, x_2, \dots, x_t)}$$

- In practice, it often makes more sense to learn (or approximate) this conditional distribution directly

Language Models

1. Convert raw text into a sequence of words (or “tokens”)

$$\mathbf{x} = [x_1, x_2, \dots, x_t]$$

2. Learn (or approximate) a joint probability distribution over sequences

$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_t)$$

- Use the chain rule of probability: predict the next word based on the previous words in the sequence

$$p(\mathbf{x}) = p(x_1)$$

$$* p(x_2|x_1)$$

$$* p(x_3|x_2, x_1)$$

$$\vdots$$

$$* p(x_t|x_{t-1}, \dots, x_1)$$

Language Models

1. Convert raw text into a sequence of words (or “tokens”)

$$\mathbf{x} = [x_1, x_2, \dots, x_t]$$

2. Learn (or approximate) a joint probability distribution over sequences

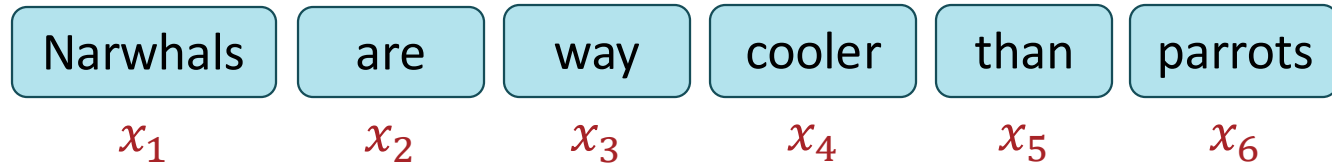
$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_t)$$

- Use the chain rule of probability: predict the next word based on the previous words in the sequence

$$p(\mathbf{x}) = p(x_1) \prod_{i=2}^t p(x_i | x_{i-1}, \dots, x_1)$$

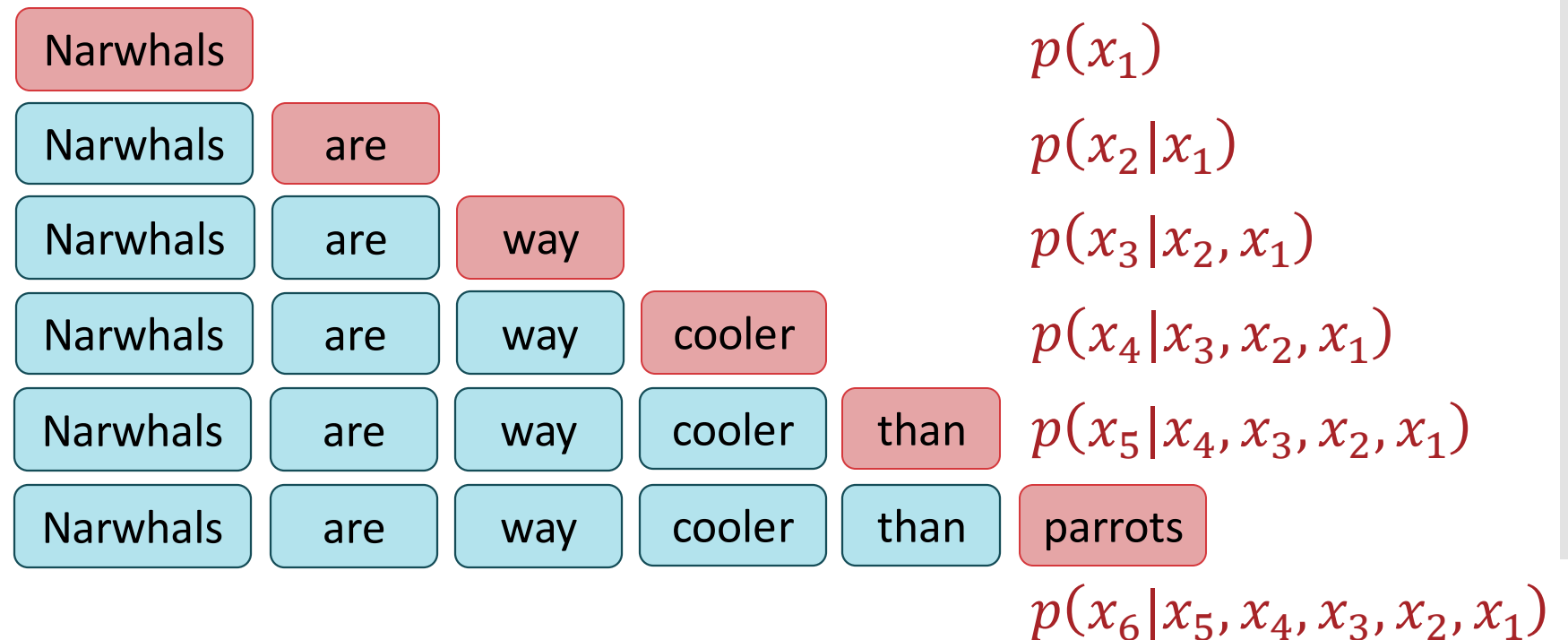
Language Models: Example

1. Convert raw text into a sequence of words (or “tokens”)



2. Learn (or approximate) a joint probability distribution over sequences

$$p(x_1, x_2, x_3, x_4, x_5, x_6) =$$



n -Gram Language Models

1. Convert raw text into a sequence of words (or “tokens”)

$$\mathbf{x} = [x_1, x_2, \dots, x_t]$$

2. Idea: approximate the chain rule using *conditional independence assumptions*

bigram model ($n = 2$): $p(\mathbf{x}) \approx p(x_1) \prod_{i=2}^t p(x_i | x_{i-1})$

n -Gram Language Models

1. Convert raw text into a sequence of words (or “tokens”)

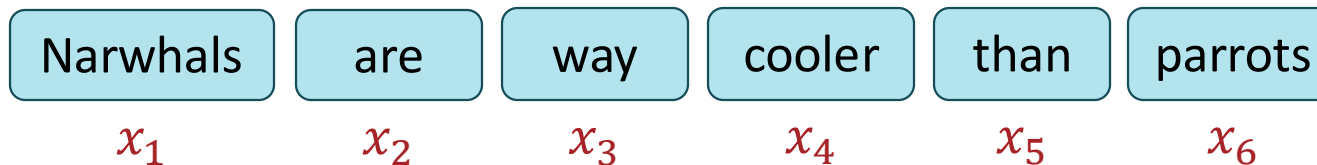
$$\mathbf{x} = [x_1, x_2, \dots, x_t]$$

2. Idea: approximate the chain rule using *conditional independence* assumptions

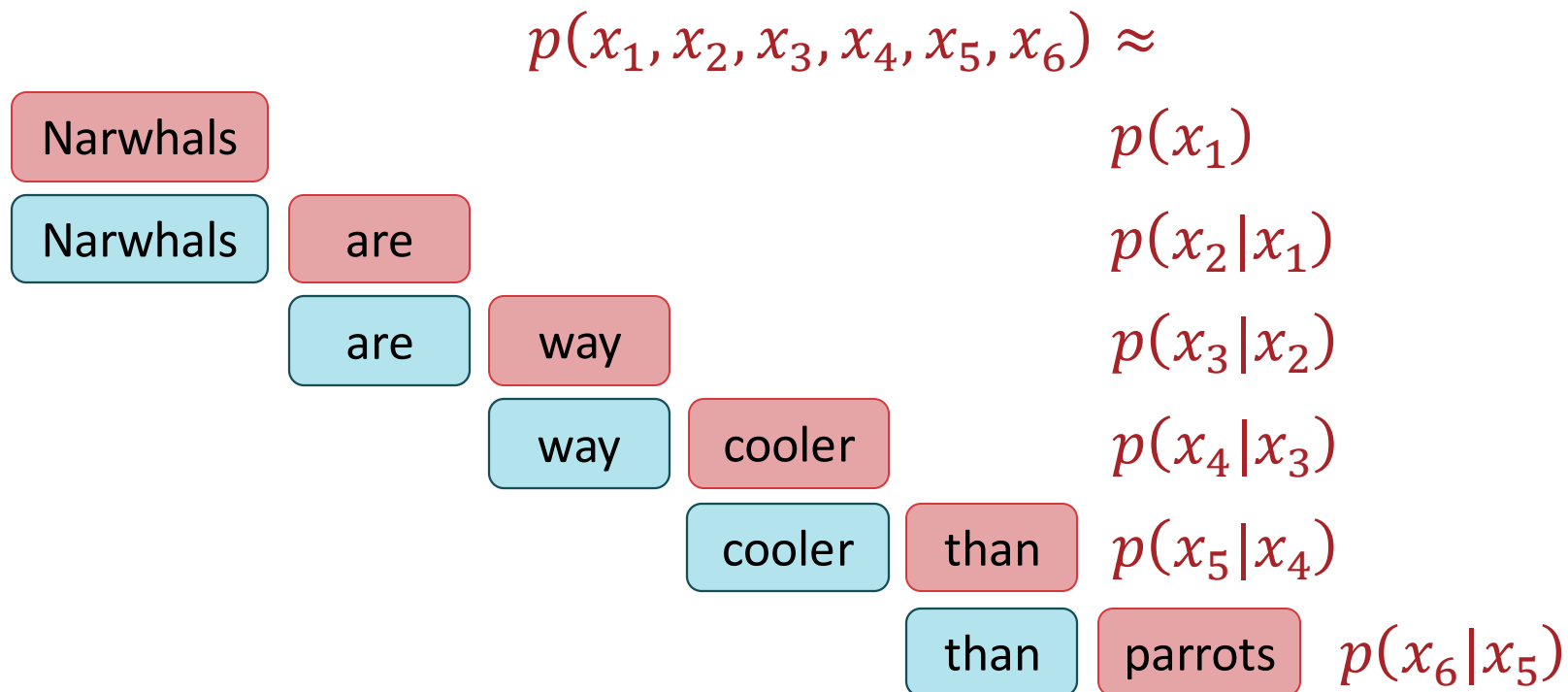
trigram model: $p(\mathbf{x}) \approx p(x_1)p(x_2|x_1) \prod_{i=3}^t p(x_i|x_{i-1}, x_{i-2})$

Bigram Language Models: Example

1. Convert raw text into a sequence of words (or “tokens”)

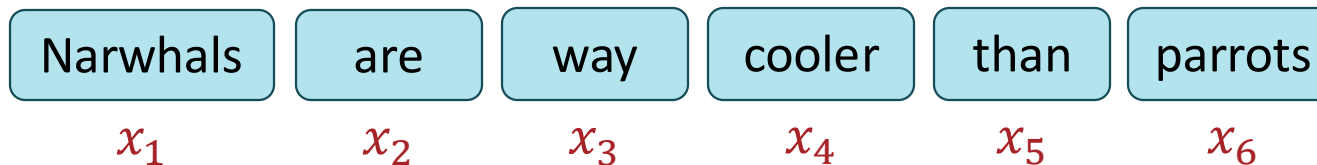


2. Idea: approximate the chain rule using *conditional independence* assumptions

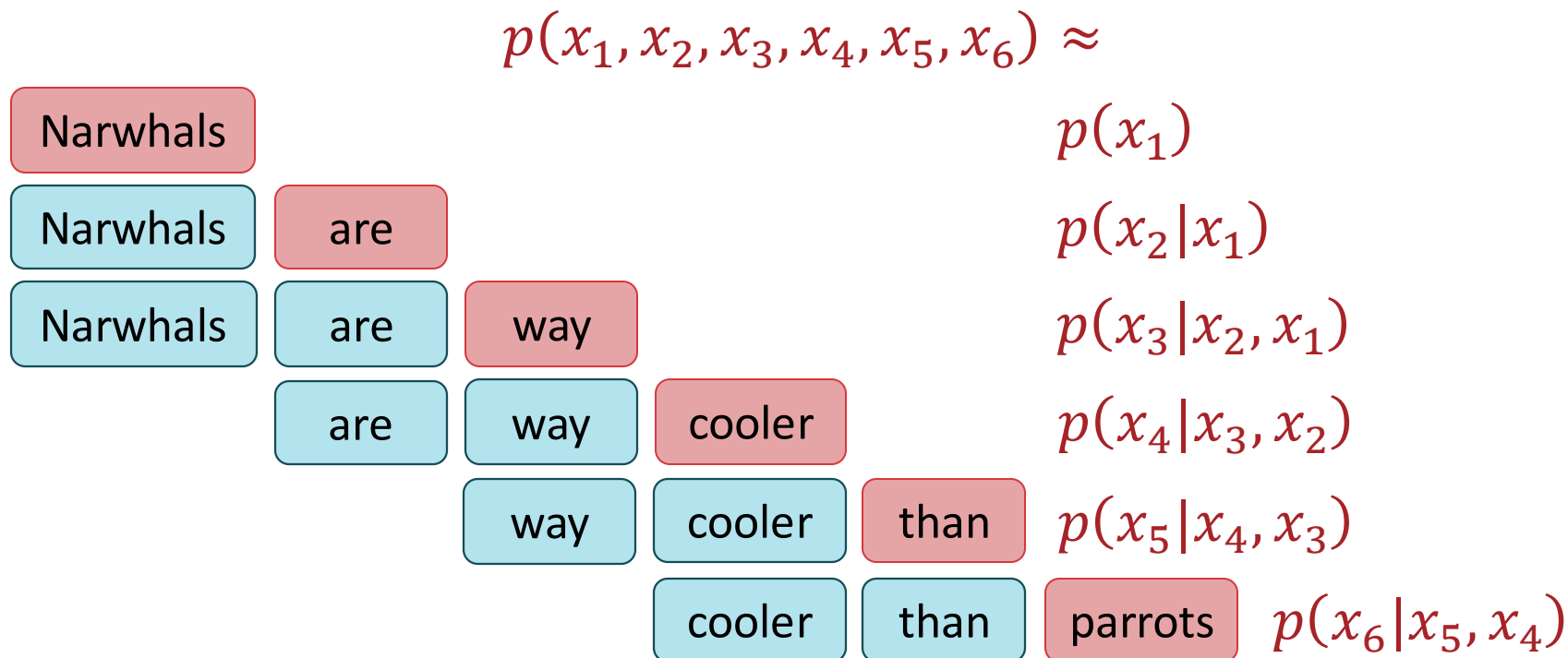


Trigram Language Models: Example

1. Convert raw text into a sequence of words (or “tokens”)



2. Idea: approximate the chain rule using *conditional independence* assumptions



n -gram Language Models: Training

- What exactly are we learning when we train an n -gram language model?
- A *lot* of probability distributions...

x_t	$p(x_t \text{narwhals, are})$		
aquatic	0.04		
awesome	x_t	$p(x_t \text{cooler, than})$	
⋮	ice	0.01	
fish	cool	x_t	$p(x_t \text{axolotls, eat})$
	⋮	shrimp	0.1
	Henry	narwhals	0.00005
		⋮	⋮
		Henry	0.0001



n -gram Language Models: Training

- How do we train an n -gram language model?
- Using training data! Simply count frequency of next words

Narwhals are big aquatic mammals that...

Who knows what **narwhals are** hiding?

Watch out, the **narwhals are** coming!

These **narwhals are** friendly!

Narwhals are a surprisingly large p...

The **narwhals are a** punk rock band

Narwhals are big fans of machine

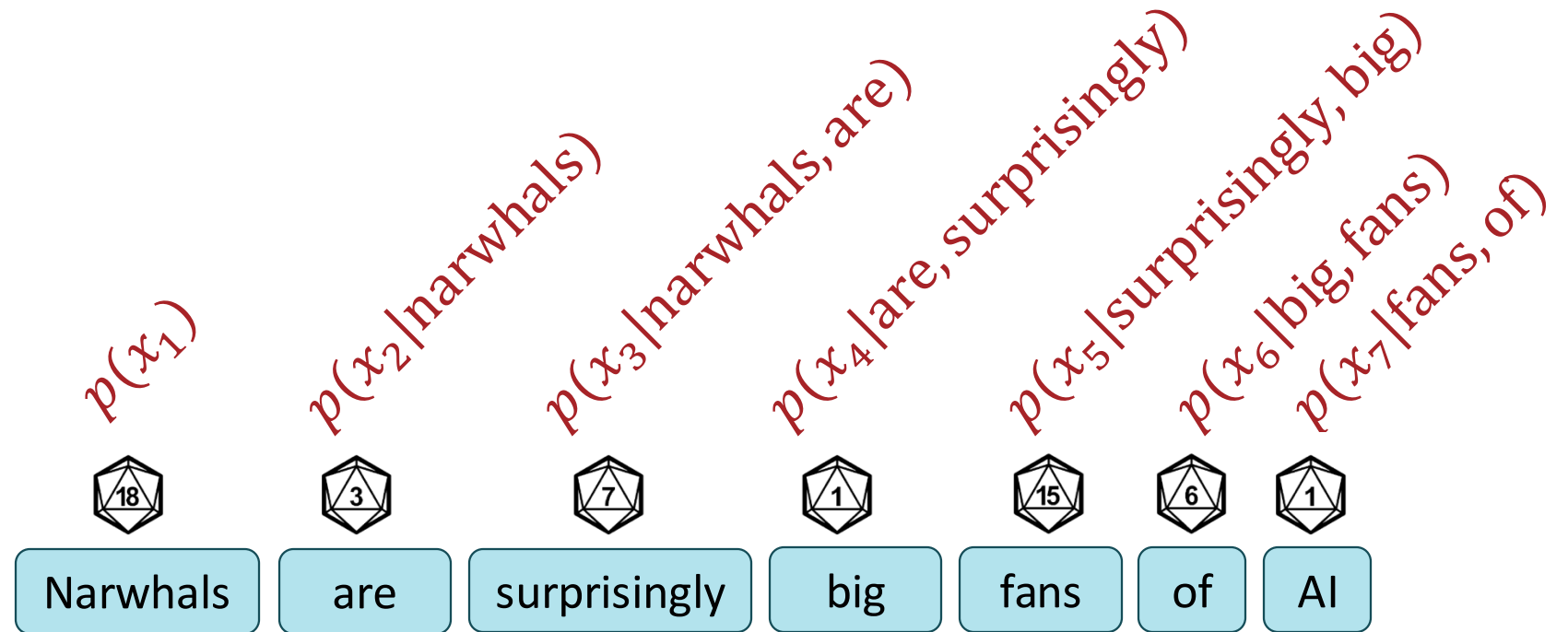
Narwhals are generated by AI.

- Why are we counting?

x_t	$p(x_t \text{narwhals, are})$
big	2/8
hiding	1/8
coming	1/8
friendly	1/8
a	2/8
generated	1/8

n -gram Language Models: Generation

- How do we generate new sequences using an n -gram language model?
- Sample from the learned distributions and use the sampled words to determine what distribution to use next



n -gram Language Models: Generation?

- How do we generate new sequences using an n -gram language model?
- Sample from the learned distributions and use the sampled words to determine what distribution to use next

Training Data (Shakespeare)

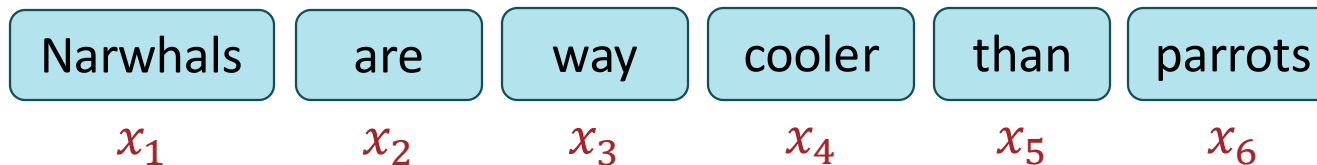
I tell you, friends, most charitable care
ave the patricians of you. For your
wants, Your suffering in this dearth,
you may as well Strike at the heaven
with your staves as lift them Against
the Roman state, whose course will on
The way it takes, cracking ten
thousand curbs Of more strong link
asunder than can ever Appear in your
impediment.

5-Gram Model

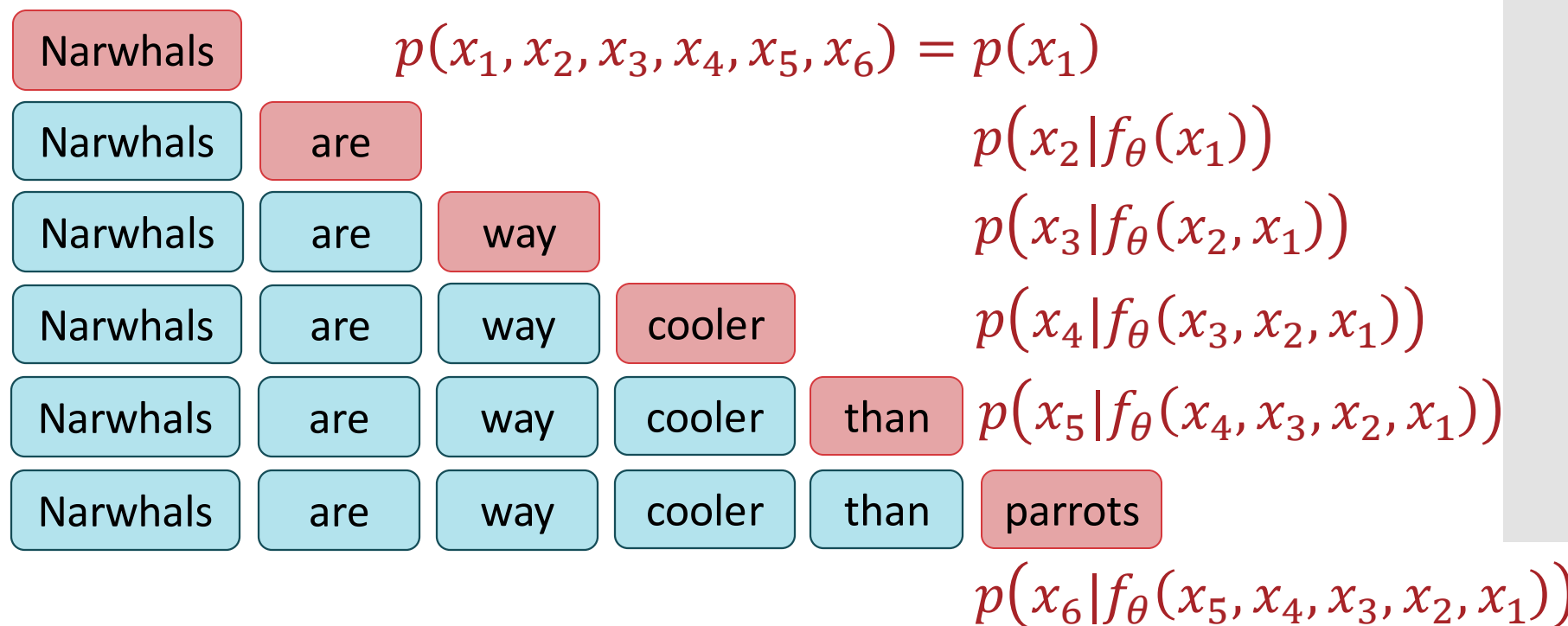
Approacheth, deny. dungy
Thither! Julius think:
grant, --0 Yead linens,
sheep's Ancient, Agreed:
Petrarch plaguy Resolved
pear! observingly honourest
adulteries wherever scabbard
guess; affirmation--his
monsieur; died. jealousy,
chequins me.

Language Models

1. Convert raw text into a sequence of words (or “tokens”)



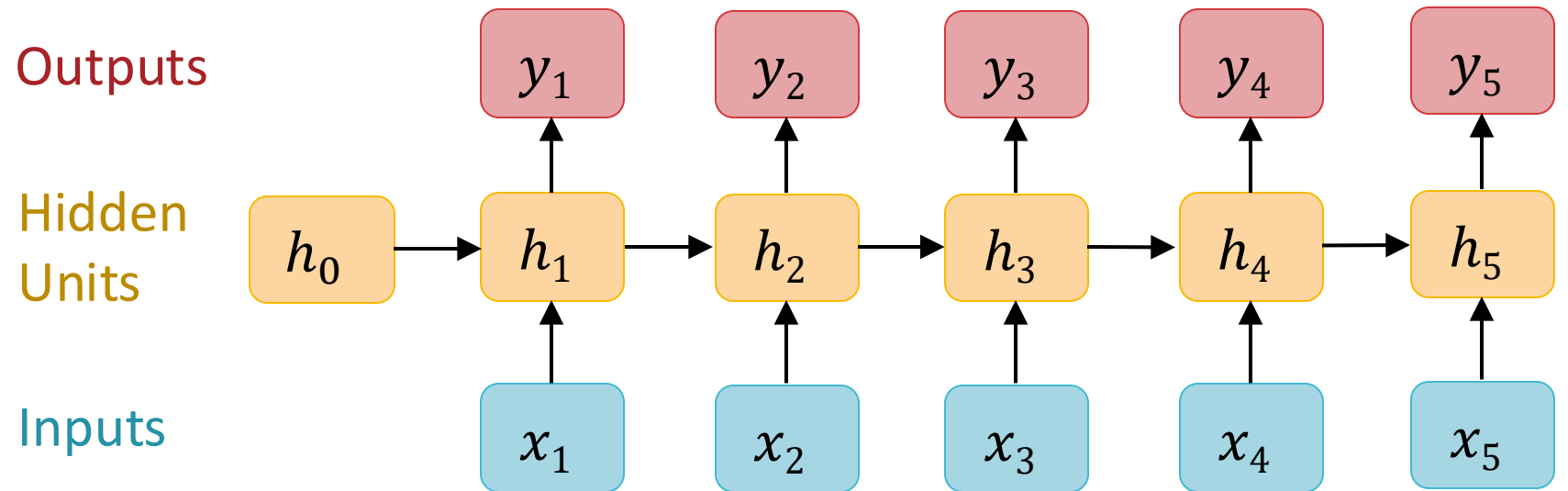
2. Idea: approximate the chain rule using a *parametric function* to compute *fixed-length* representations of prior words and condition on those to compute next words



Recurrent Neural Networks

$$h_t = \phi(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$

$$y_t = \psi(W_{hy}h_t + b_y)$$



Recurrent Neural Network Language Models

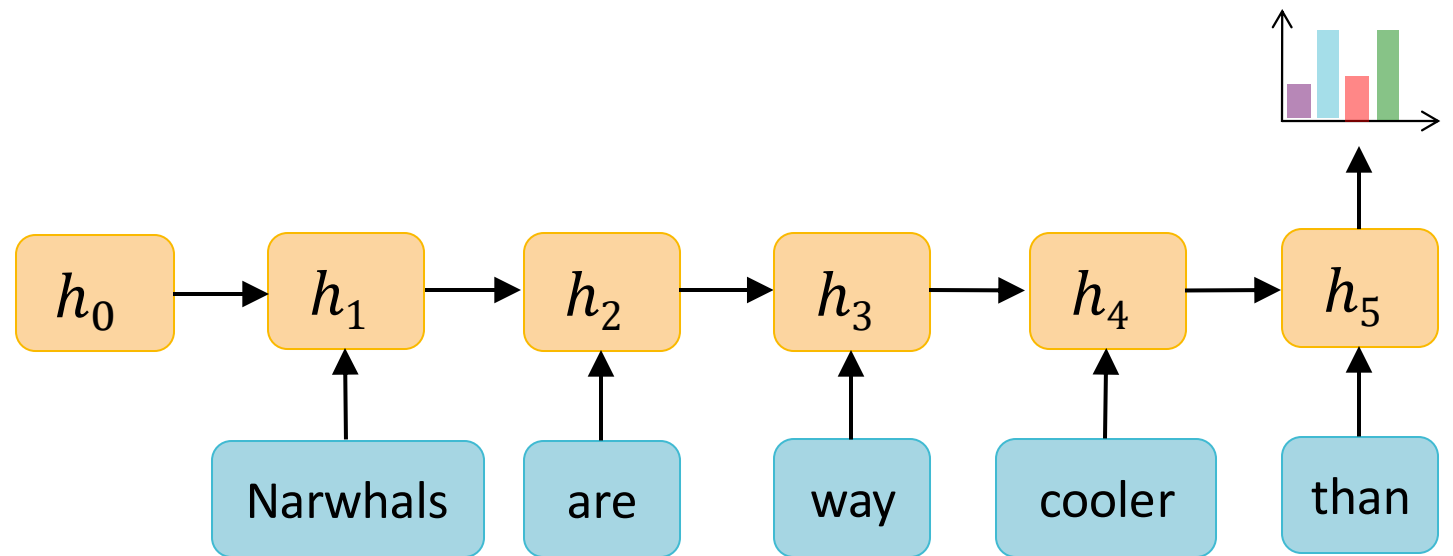
- How can we get a distribution $p(x_t|h_{t-1} = f_{\theta}(x_1, \dots, x_{t-1}))$ from some fixed-length vector h_{t-1} ?

$$p(x_6|h_5 = f_{\theta}(x_1, \dots, x_5))$$

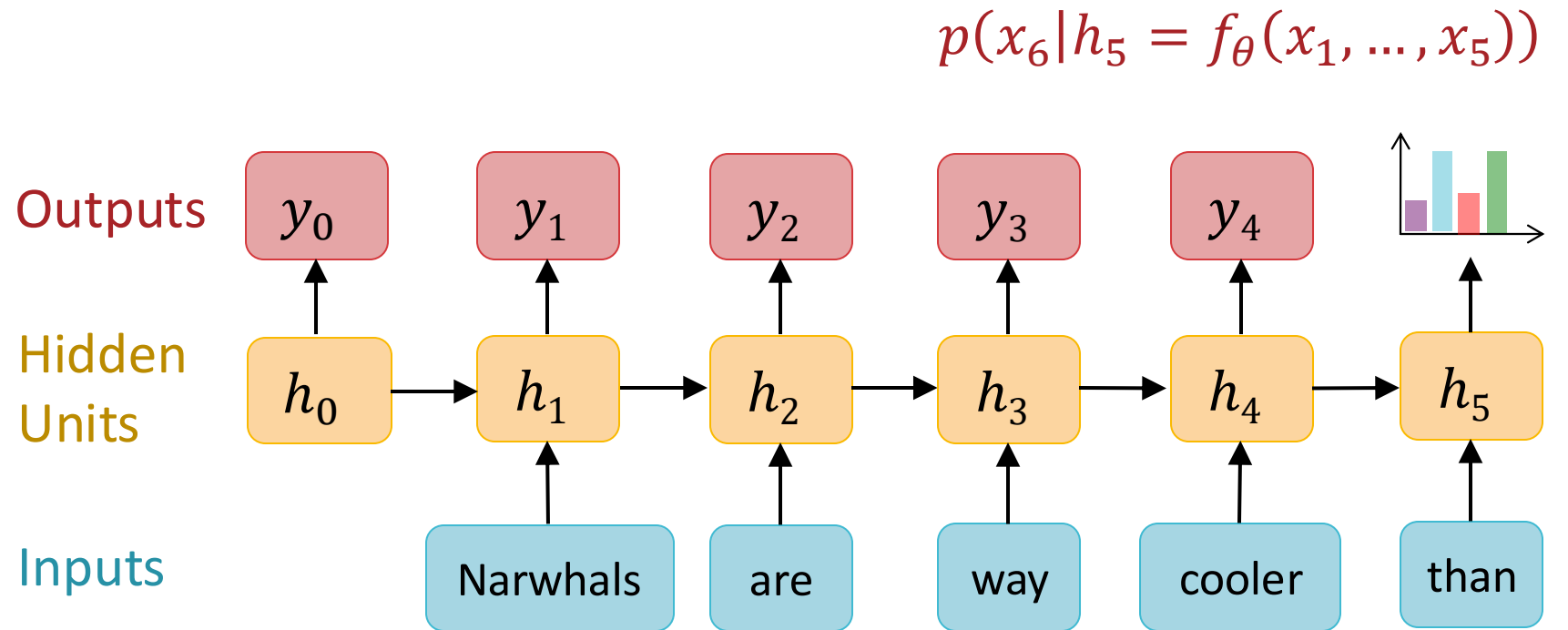
Outputs

Hidden
Units

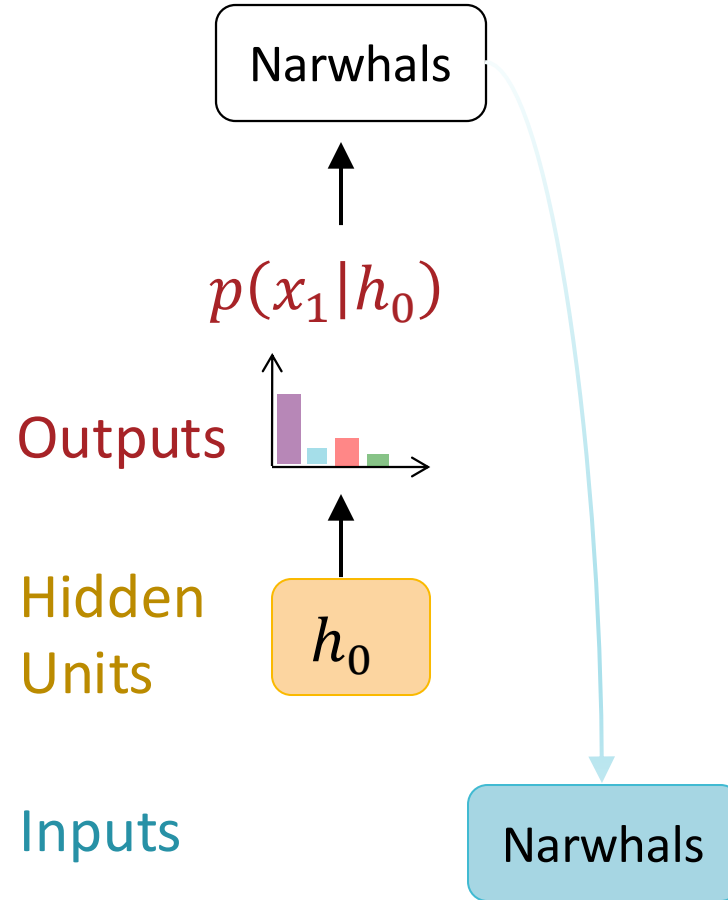
Inputs



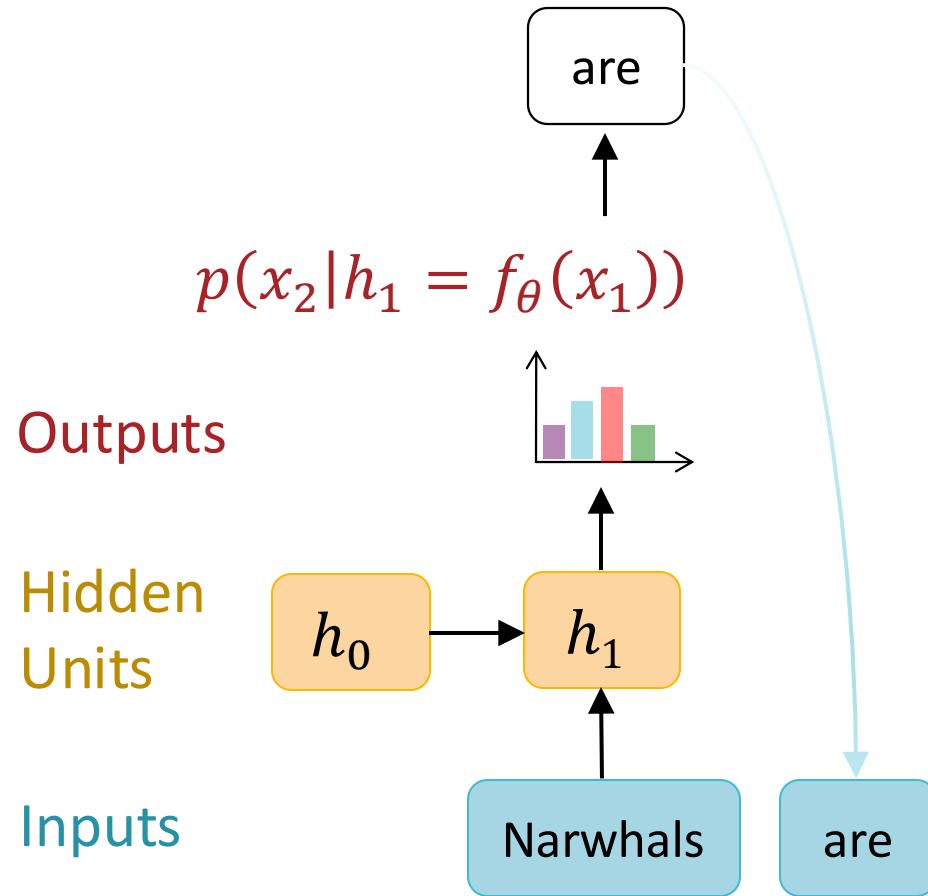
Recurrent Neural Network Language Models



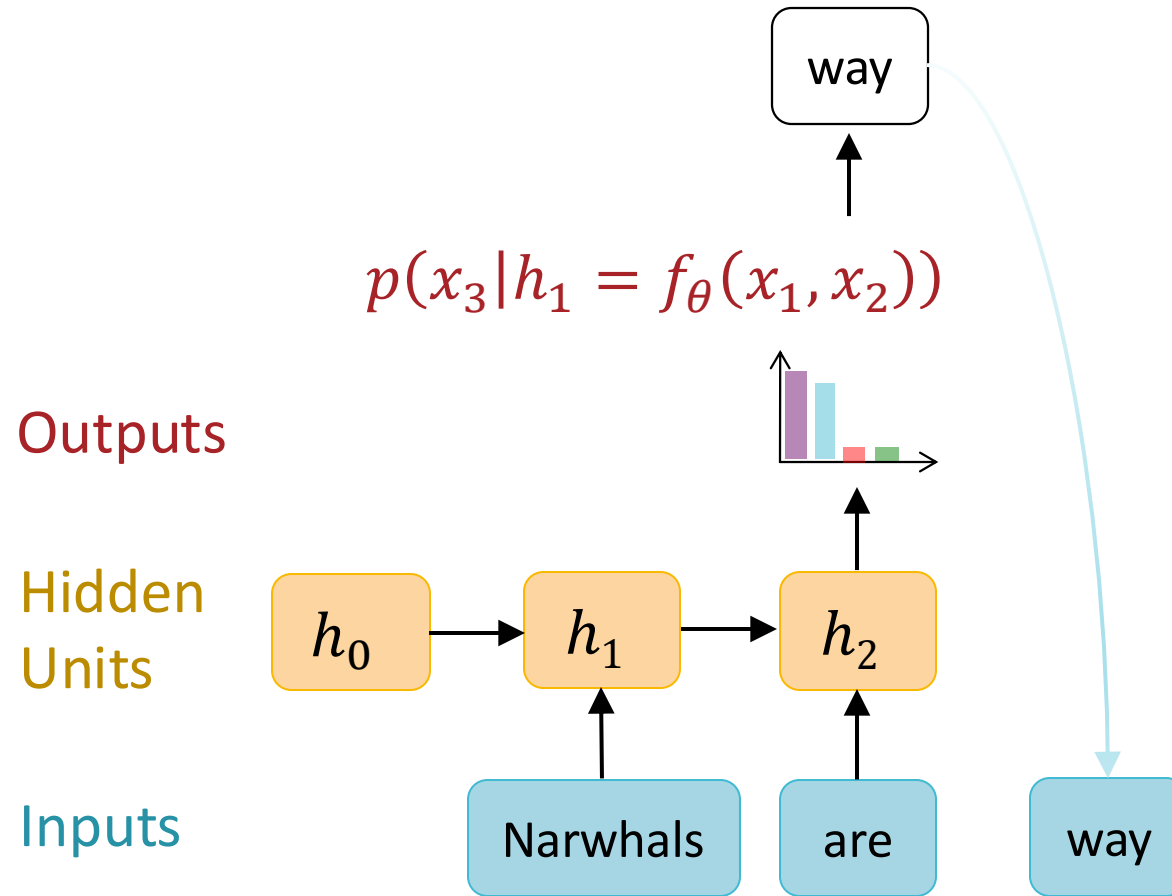
Recurrent Neural Network Language Models: Generation



Recurrent Neural Network Language Models: Generation

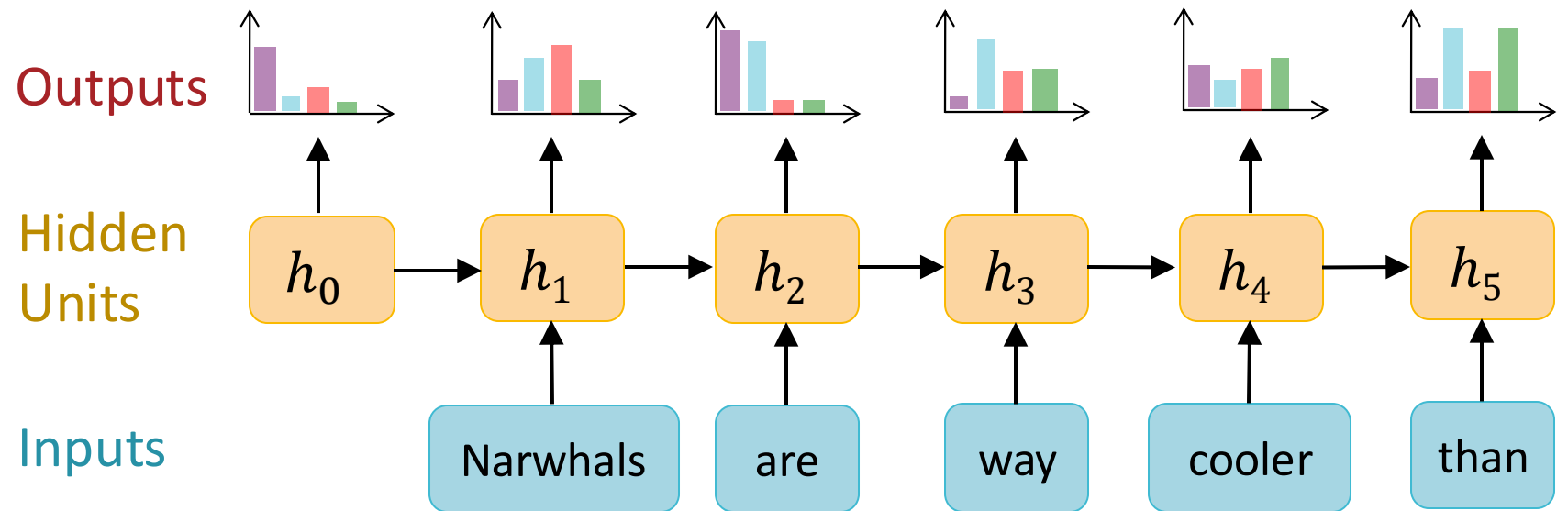


Recurrent Neural Network Language Models: Generation



Recurrent Neural Network Language Models: Generation

- How do we generate new sequences using an RNN language model?
- Exactly the same way we did for an n -gram language model, by sampling from some learned probability distributions over next words!



Shakespeare's As You Like It

VIOLA: Why, Salisbury must find his flesh and thought
That which I am not apt, not a man and in fire, To show
the reining of the raven and the wars To grace my hand
reproach within, and not a fair are hand, That Caesar
and my goodly father's world; When I was heaven of
presence and our fleets, We spare with hours, but cut
thy council I am great, Murdered and by thy master's
ready there My power to give thee but so much as hell:
Some service in the noble bondman here, Would show
him to her wine.

KING LEAR: O, if you were a feeble sight, the courtesy
of your law, Your sight and several breath, will wear the
gods With his heads, and my hands are wonder'd at
the deeds, So drop upon your lordship's head, and your
opinion Shall be against your honour.

RNN-LM Sample

CHARLES: Marry, do I, sir; and I came to acquaint you
with a matter. I am given, sir, secretly to understand
that your younger brother Orlando hath a disposition
to come in disguised against me to try a fall. To-
morrow, sir, I wrestle for my credit; and he that escapes
me without some broken limb shall acquit him well.
Your brother is but young and tender; and, for your
love, I would be loath to foil him, as I must, for my own
honour, if he come in: therefore, out of my love to you,
I came hither to acquaint you withal, that either you
might stay him from his intendment or brook such
disgrace well as he shall run into, in that it is a thing of
his own search and altogether against my will.

TOUCHSTONE: For my part, I had rather bear with you
than bear you; yet I should bear no cross if I did bear
you, for I think you have no money in your purse.

RNN Language Models: Generation

RNN-LM Sample?

VIOLA: Why, Salisbury must find his flesh and thought
That which I am not apt, not a man and in fire, To show
the reining of the raven and the wars To grace my hand
reproach within, and not a fair are hand, That Caesar
and my goodly father's world; When I was heaven of
presence and our fleets, We spare with hours, but cut
thy council I am great, Murdered and by thy master's
ready there My power to give thee but so much as hell:
Some service in the noble bondman here, Would show
him to her wine.

KING LEAR: O, if you were a feeble sight, the courtesy
of your law, Your sight and several breath, will wear the
gods With his heads, and my hands are wonder'd at
the deeds, So drop upon your lordship's head, and your
opinion Shall be against your honour.

Shakespeare's As You Like It?

CHARLES: Marry, do I, sir; and I came to acquaint you
with a matter. I am given, sir, secretly to understand
that your younger brother Orlando hath a disposition
to come in disguised against me to try a fall. To-
morrow, sir, I wrestle for my credit; and he that escapes
me without some broken limb shall acquit him well.
Your brother is but young and tender; and, for your
love, I would be loath to foil him, as I must, for my own
honour, if he come in: therefore, out of my love to you,
I came hither to acquaint you withal, that either you
might stay him from his intendment or brook such
disgrace well as he shall run into, in that it is a thing of
his own search and altogether against my will.

TOUCHSTONE: For my part, I had rather bear with you
than bear you; yet I should bear no cross if I did bear
you, for I think you have no money in your purse.

RNN Language Models: Generation

RNN-LM Sample?

VIOLA: Why, Salisbury must find his flesh and thought
That which I am not apt, not a man and in fire, To show
the reining of the raven and the wars To grace my hand
reproach within, and not a fair are hand, That Caesar
and my goodly father's world; When I was heaven of
presence and our fleets, We spare with hours, but cut
thy council I am great, Murdered and by thy master's
ready there My power to give thee but so much as hell:
Some service in the noble bondman here, Would show
him to her wine.

KING LEAR: O, if you were a feeble sight, the courtesy
of your law, Your sight and several breath, will wear the
gods With his heads, and my hands are wonder'd at
the deeds, So drop upon your lordship's head, and your
opinion Shall be against your honour.

Shakespeare's As You Like It?

CHARLES: Marry, do I, sir; and I came to acquaint you
with a matter. I am given, sir, secretly to understand
that your younger brother Orlando hath a disposition
to come in disguised against me to try a fall. To-
morrow, sir, I wrestle for my credit; and he that escapes
me without some broken limb shall acquit him well.
Your brother is but young and tender; and, for your
love, I would be loath to foil him, as I must, for my own
honour, if he come in: therefore, out of my love to you,
I came hither to acquaint you withal, that either you
might stay him from his intendment or brook such
disgrace well as he shall run into, in that it is a thing of
his own search and altogether against my will.

TOUCHSTONE: For my part, I had rather bear with you
than bear you; yet I should bear no cross if I did bear
you, for I think you have no money in your purse.

RNN Language Models: Training?

Backpropagation

Forward Computation

1. Write an **algorithm** for evaluating the function $y = f(\mathbf{x})$.
 - The algorithm defines the **computation graph**, a **directed acyclic graph**, where each variable is a node
2. Visit each node in **topological order**:
For variable u_i with inputs v_1, \dots, v_N compute
 $u_i = g_i(v_1, \dots, v_N)$ and store the result at that node

Backward Computation (Version A)

1. **Initialize** $dy/dy = 1$.
2. Visit each node v_j in **reverse topological order**:
 - Let u_1, \dots, u_M denote all the nodes with v_j as an input
Assume $y = h(\mathbf{u}) = h(u_1, \dots, u_M)$ and
 $u_i = g_i(\mathbf{v}) = g_i(v_1, \dots, v_j, \dots, v_N) \forall i$
3. Compute $\frac{dy}{dv_j} = \sum_{i=1}^M \frac{dy}{du_i} \frac{du_i}{dv_j}$
 - We already know $dy/du_i \forall i$ and du_i/dv_j is (typically) easy to compute
4. **Return** partial derivatives dy/du_i for all variables

Backpropagation

Forward Computation

1. Write an **algorithm** for evaluating the function $y = f(\mathbf{x})$.
 - The algorithm defines the **computation graph**, a **directed acyclic graph**, where each variable is a node
2. Visit each node in **topological order**:
For variable u_i with inputs v_1, \dots, v_N compute $u_i = g_i(v_1, \dots, v_N)$ and store the result at that node

Backward Computation (Version B)

1. **Initialize** $dy/dy = 1$.
2. Visit each node v_j in **reverse topological order**
3. For each variable $u_i = g(\mathbf{v}) = g_i(v_1, \dots, v_j, \dots, v_N)$, increment dy/dv_j by $(dy/du_i)(du_i/dv_j)$
 - a) We already know $dy/du_i \forall i$ and du_i/dv_j is (typically) easy to compute
4. **Return** partial derivatives dy/du_i for all variables

- Key takeaway: backpropagation makes efficient reuse of intermediate quantities when computing the necessary gradients to train a neural network

Backpropagation: Procedural Method

Algorithm 1 Forward Computation

```
1: procedure NNFORWARD(Training example  $(\mathbf{x}, \mathbf{y})$ , Params  $\alpha, \beta$ )
2:    $\mathbf{a} = \alpha \mathbf{x}$ 
3:    $\mathbf{z} = \sigma(\mathbf{a})$ 
4:    $\mathbf{b} = \beta \mathbf{z}$ 
5:    $\hat{\mathbf{y}} = \text{softmax}(\mathbf{b})$ 
6:    $J = -\mathbf{y}^T \log \hat{\mathbf{y}}$ 
7:    $\mathbf{o} = \text{object}(\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J)$ 
8:   return intermediate quantities  $\mathbf{o}$ 
```

Algorithm 2 Backpropagation

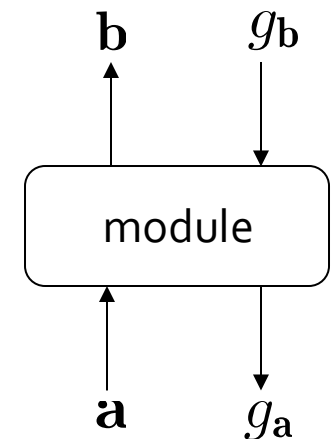
```
1: procedure NNBACKWARD(Training example  $(\mathbf{x}, \mathbf{y})$ , Params  $\alpha, \beta$ ,  
Intermediates  $\mathbf{o}$ )
2:   Place intermediate quantities  $\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J$  in  $\mathbf{o}$  in scope
3:    $\mathbf{g}_{\hat{\mathbf{y}}} = -\mathbf{y} \div \hat{\mathbf{y}}$ 
4:    $\mathbf{g}_{\mathbf{b}} = \mathbf{g}_{\hat{\mathbf{y}}}^T (\text{diag}(\hat{\mathbf{y}}) - \hat{\mathbf{y}}\hat{\mathbf{y}}^T)$ 
5:    $\mathbf{g}_{\beta} = \mathbf{g}_{\mathbf{b}}^T \mathbf{z}^T$ 
6:    $\mathbf{g}_{\mathbf{z}} = \beta^T \mathbf{g}_{\mathbf{b}}^T$ 
7:    $\mathbf{g}_{\mathbf{a}} = \mathbf{g}_{\mathbf{z}} \odot \mathbf{z} \odot (1 - \mathbf{z})$ 
8:    $\mathbf{g}_{\alpha} = \mathbf{g}_{\mathbf{a}} \mathbf{x}^T$ 
9:   return parameter gradients  $\mathbf{g}_{\alpha}, \mathbf{g}_{\beta}$ 
```

Issues:

1. Hard to reuse / adapt for other models
2. Hard to optimize individual steps
3. Hard to debug using the finite-difference check

Module-based AutoDiff

- Key Idea:
 - componentize the computation of the neural-network into layers
 - each layer consolidates multiple **real-valued nodes** in the computation graph (a subset of them) into one **vector-valued node** (aka. a **module**)
- Each **module** is capable of two actions:
 - **Forward computation of the output given some input**
 - **Backward computation of the gradient with respect to the input given the gradient with respect to the output**



Module-based AutoDiff

Linear Module The linear layer has two inputs: a vector \mathbf{a} and parameters $\omega \in \mathbb{R}^{B \times A}$. The output \mathbf{b} is not used by LINEARBACKWARD, but we pass it in for consistency of form.

```
1: procedure LINEARFORWARD( $\mathbf{a}$ ,  $\omega$ )
2:    $\mathbf{b} = \omega \mathbf{a}$ 
3:   return  $\mathbf{b}$ 
4: procedure LINEARBACKWARD( $\mathbf{a}$ ,  $\omega$ ,  $\mathbf{b}$ ,  $\mathbf{g}_b$ )
5:    $\mathbf{g}_\omega = \mathbf{g}_b \mathbf{a}^T$ 
6:    $\mathbf{g}_a = \omega^T \mathbf{g}_b$ 
7:   return  $\mathbf{g}_\omega$ ,  $\mathbf{g}_a$ 
```

Softmax Module The softmax layer has only one input vector \mathbf{a} . For any vector $\mathbf{v} \in \mathbb{R}^D$, we have that $\text{diag}(\mathbf{v})$ returns a $D \times D$ diagonal matrix whose diagonal entries are v_1, v_2, \dots, v_D and whose non-diagonal entries are zero.

```
1: procedure SOFTMAXFORWARD( $\mathbf{a}$ )
2:    $\mathbf{b} = \text{softmax}(\mathbf{a})$ 
3:   return  $\mathbf{b}$ 
4: procedure SOFTMAXBACKWARD( $\mathbf{a}$ ,  $\mathbf{b}$ ,  $\mathbf{g}_b$ )
5:    $\mathbf{g}_a = \mathbf{g}_b^T (\text{diag}(\mathbf{b}) - \mathbf{b}\mathbf{b}^T)$ 
6:   return  $\mathbf{g}_a$ 
```

Sigmoid Module The sigmoid layer has only one input vector \mathbf{a} . Below σ is the sigmoid applied element-wise, and \odot is element-wise multiplication s.t. $\mathbf{u} \odot \mathbf{v} = [u_1 v_1, \dots, u_M v_M]$.

```
1: procedure SIGMOIDFORWARD( $\mathbf{a}$ )
2:    $\mathbf{b} = \sigma(\mathbf{a})$ 
3:   return  $\mathbf{b}$ 
4: procedure SIGMOIDBACKWARD( $\mathbf{a}$ ,  $\mathbf{b}$ ,  $\mathbf{g}_b$ )
5:    $\mathbf{g}_a = \mathbf{g}_b \odot \mathbf{b} \odot (1 - \mathbf{b})$ 
6:   return  $\mathbf{g}_a$ 
```

Cross-Entropy Module The cross-entropy layer has two inputs: a gold one-hot vector \mathbf{a} and a predicted probability distribution $\hat{\mathbf{a}}$. Its output $b \in \mathbb{R}$ is a scalar. Below \div is element-wise division. The output b is not used by CROSSENTROPYBACKWARD, but we pass it in for consistency of form.

```
1: procedure CROSSENTROPYFORWARD( $\mathbf{a}$ ,  $\hat{\mathbf{a}}$ )
2:    $b = -\mathbf{a}^T \log \hat{\mathbf{a}}$ 
3:   return  $\mathbf{b}$ 
4: procedure CROSSENTROPYBACKWARD( $\mathbf{a}$ ,  $\hat{\mathbf{a}}$ ,  $b$ ,  $\mathbf{g}_b$ )
5:    $\mathbf{g}_{\hat{\mathbf{a}}} = -\mathbf{g}_b (\mathbf{a} \div \hat{\mathbf{a}})$ 
6:   return  $\mathbf{g}_a$ 
```

Module-based AutoDiff

Algorithm 1 Forward Computation

```
1: procedure NNFORWARD(Training example  $(\mathbf{x}, \mathbf{y})$ , Parameters  $\alpha, \beta$ )
2:    $\mathbf{a} = \text{LINEARFORWARD}(\mathbf{x}, \alpha)$ 
3:    $\mathbf{z} = \text{SIGMOIDFORWARD}(\mathbf{a})$ 
4:    $\mathbf{b} = \text{LINEARFORWARD}(\mathbf{z}, \beta)$ 
5:    $\hat{\mathbf{y}} = \text{SOFTMAXFORWARD}(\mathbf{b})$ 
6:    $J = \text{CROSSENTROPYFORWARD}(\mathbf{y}, \hat{\mathbf{y}})$ 
7:    $\mathbf{o} = \text{object}(\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J)$ 
8:   return intermediate quantities  $\mathbf{o}$ 
```

Algorithm 2 Backpropagation

```
1: procedure NNBACKWARD(Training example  $(\mathbf{x}, \mathbf{y})$ , Parameters  $\alpha, \beta$ , Intermediates  $\mathbf{o}$ )
2:   Place intermediate quantities  $\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J$  in  $\mathbf{o}$  in scope
3:    $g_J = \frac{dJ}{dJ} = 1$  ▷ Base case
4:    $\mathbf{g}_{\hat{\mathbf{y}}} = \text{CROSSENTROPYBACKWARD}(\mathbf{y}, \hat{\mathbf{y}}, J, g_J)$ 
5:    $\mathbf{g}_{\mathbf{b}} = \text{SOFTMAXBACKWARD}(\mathbf{b}, \hat{\mathbf{y}}, \mathbf{g}_{\hat{\mathbf{y}}})$ 
6:    $\mathbf{g}_{\beta}, \mathbf{g}_{\mathbf{z}} = \text{LINEARBACKWARD}(\mathbf{z}, \mathbf{b}, \mathbf{g}_{\mathbf{b}})$ 
7:    $\mathbf{g}_{\mathbf{a}} = \text{SIGMOIDBACKWARD}(\mathbf{a}, \mathbf{z}, \mathbf{g}_{\mathbf{z}})$ 
8:    $\mathbf{g}_{\alpha}, \mathbf{g}_{\mathbf{x}} = \text{LINEARBACKWARD}(\mathbf{x}, \mathbf{a}, \mathbf{g}_{\mathbf{a}})$  ▷ We discard  $\mathbf{g}_{\mathbf{x}}$ 
9:   return parameter gradients  $\mathbf{g}_{\alpha}, \mathbf{g}_{\beta}$ 
```

1. Easy to reuse / adapt for other models
2. Individual layers are easier to optimize
3. Simple to debug: just run a finite-difference check on each layer separately

Module-based AutoDiff (OOP Version)

Object-Oriented Implementation:

- Let each module be an object and allow the **control flow** of the program to define the computation graph
- **No longer need to implement NNBackward(\cdot)**, just follow the computation graph in reverse topological order

```
1 class Sigmoid(Module)
2   method forward(a)
3      $\mathbf{b} = \sigma(\mathbf{a})$ 
4     return  $\mathbf{b}$ 
5   method backward(a, b, gb)
6      $\mathbf{g}_a = \mathbf{g}_b \odot \mathbf{b} \odot (1 - \mathbf{b})$ 
7     return  $\mathbf{g}_a$ 
```

```
1 class Softmax(Module)
2   method forward(a)
3      $\mathbf{b} = \text{softmax}(\mathbf{a})$ 
4     return  $\mathbf{b}$ 
5   method backward(a, b, gb)
6      $\mathbf{g}_a = \mathbf{g}_b^T (\text{diag}(\mathbf{b}) - \mathbf{b}\mathbf{b}^T)$ 
7     return  $\mathbf{g}_a$ 
```

```
1 class Linear(Module)
2   method forward(a,  $\omega$ )
3      $\mathbf{b} = \omega \mathbf{a}$ 
4     return  $\mathbf{b}$ 
5   method backward(a,  $\omega$ , b, gb)
6      $\mathbf{g}_\omega = \mathbf{g}_b \mathbf{a}^T$ 
7      $\mathbf{g}_a = \omega^T \mathbf{g}_b$ 
8     return  $\mathbf{g}_\omega, \mathbf{g}_a$ 
```

```
1 class CrossEntropy(Module)
2   method forward(a,  $\hat{\mathbf{a}}$ )
3      $\mathbf{b} = -\mathbf{a}^T \log \hat{\mathbf{a}}$ 
4     return  $\mathbf{b}$ 
5   method backward(a,  $\hat{\mathbf{a}}$ , b, gb)
6      $\mathbf{g}_{\hat{\mathbf{a}}} = -g_b(\mathbf{a} \div \hat{\mathbf{a}})$ 
7     return  $\mathbf{g}_a$ 
```

Module-based AutoDiff (OOP Version)

```
1 class NeuralNetwork(Module):
2
3     method init()
4         lin1_layer = Linear()
5         sig_layer = Sigmoid()
6         lin2_layer = Linear()
7         soft_layer = Softmax()
8         ce_layer = CrossEntropy()
9
10    method forward(Tensor x, Tensor y, Tensor  $\alpha$ , Tensor  $\beta$ )
11        a = lin1_layer.apply_fwd(x,  $\alpha$ )
12        z = sig_layer.apply_fwd(a)
13        b = lin2_layer.apply_fwd(z,  $\beta$ )
14         $\hat{y}$  = soft_layer.apply_fwd(b)
15        J = ce_layer.apply_fwd(y,  $\hat{y}$ )
16        return J.out_tensor
17
18    method backward(Tensor x, Tensor y, Tensor  $\alpha$ , Tensor  $\beta$ )
19        tape_bwd()
20        return lin1_layer.in_gradients[1], lin2_layer.in_gradients[1]
```

Module-based AutoDiff (OOP Version)

```
1 global tape = stack()
2
3 class Module:
4
5     method init()
6         out_tensor = null
7         out_gradient = 1
8
9     method apply_fwd(List in_modules)
10        in_tensors = [x.out_tensor for x in in_modules]
11        out_tensor = forward(in_tensors)
12        tape.push(self)
13        return self
14
15    method apply_bwd():
16        in_gradients = backward(in_tensors, out_tensor, out_gradient)
17        for i in 1, ..., len(in_modules):
18            in_modules[i].out_gradient += in_gradients[i]
19        return self
20
21 function tape_bwd():
22     while len(tape) > 0
23         m = tape.pop()
24         m.apply_bwd()
```

PyTorch

```
class NeuralNetwork(nn.Module):
    def __init__(self):
        super().__init__()
        self.flatten = nn.Flatten()
        self.linear1 = nn.Linear(28*28, 512)
        self.sigmoid = nn.Sigmoid()
        self.linear2 = nn.Linear(512, 512)

    def forward(self, x):
        x = self.flatten(x)
        a = self.linear1(x)
        z = self.sigmoid(a)
        b = self.linear2(z)
        return b

def one_step_of_sgd(X, y):
    model = NeuralNetwork()
    loss_fn = nn.CrossEntropyLoss()
    optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)

    pred = model(X)
    loss = loss_fn(pred, y)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

PyTorch: FAQ

- Q: Why don't we call `linear.forward()` in PyTorch?
- A: There's a special method in Python `__call__` that allows you to define what happens when you treat an object as if it were a function so `linear(x)` is equivalent to `linear.__call__(x)` and PyTorch defines every Module's `__call__` method to be something like

```
def __call__(self):  
    self.forward()
```

PyTorch: FAQ

- Q: Why don't we pass in the parameters to a Module?
- A: This is just a design choice; in PyTorch, you store the parameters inside the Module and “mark” them as parameters that should contribute to the eventual gradient used by an optimizer