



Edge AI Workshop

Lecture II - Language Models

Adrian Cosma





Workshop Overview

- 1. Part 0: Introduction to the Team & Rust for Edge AI**
- 2. Part I: Lecture on Computer Vision**
 - a. Main problems in Computer Vision
 - b. What exactly is a neural network? (CNNs / Transformers)
 - c. What exactly is an image embedding?
 - d. Computer vision on the Edge
- 3. Hands-On I: Air-gapped Face recognition on the Pi**
- 4. Part II: Lecture on Natural Language Processing**
 - a. A bit of history & development of modern LLMs
 - b. How does an LLM work? Tokenizers, pretraining, post-training
 - c. Context Engineering: Tool calling, RAG
 - d. Libraries: tokenizers-rs, llama.cpp
- 5. Hands-On II: Chat with a LLM on Pi**





Workshop Overview

- 1. Part 0: Introduction to the Team & Rust for Edge AI**
- 2. Part I: Lecture on Computer Vision**
 - a. Main problems in Computer Vision
 - b. What exactly is a neural network? (CNNs / Transformers)
 - c. What exactly is an image embedding?
 - d. Computer vision on the Edge
- 3. Hands-On I: Air-gapped Face recognition on the Pi**
- 4. Part II: Lecture on Natural Language Processing**
 - a. A bit of history & development of modern LLMs
 - b. How does an LLM work? Tokenizers, pretraining, post-training
 - c. Context Engineering: Tool calling, RAG
 - d. Libraries: tokenizers-rs, llama.cpp
- 5. Hands-On II: Chat with a LLM on Pi**



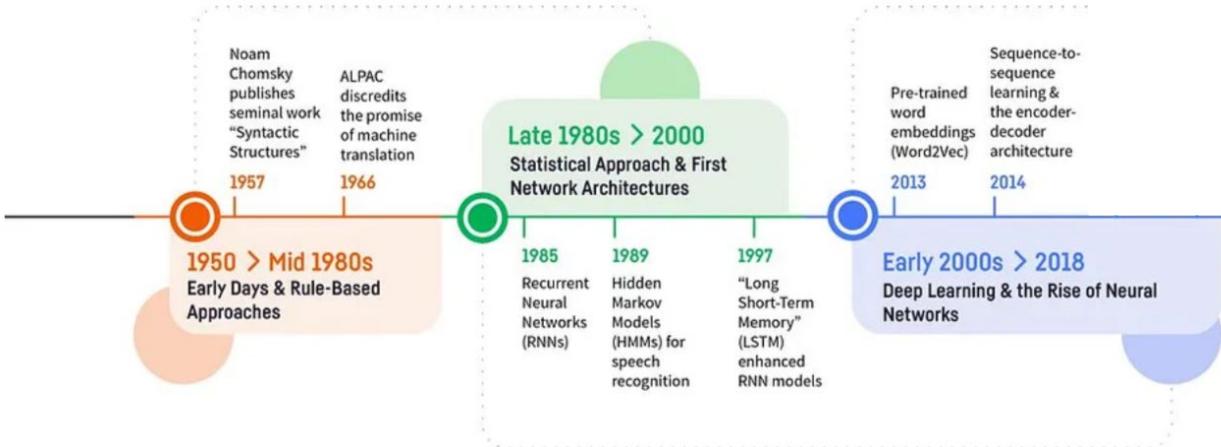


This lecture

- The historical developments of the ideas behind modern language models
 - Mostly from OpenAI, since their path was quite clear and their models were highly influential to the research community
- High level overview of how modern LLMs are trained and fine-tuned
- Prompting methods
- Patterns for deployment on the edge

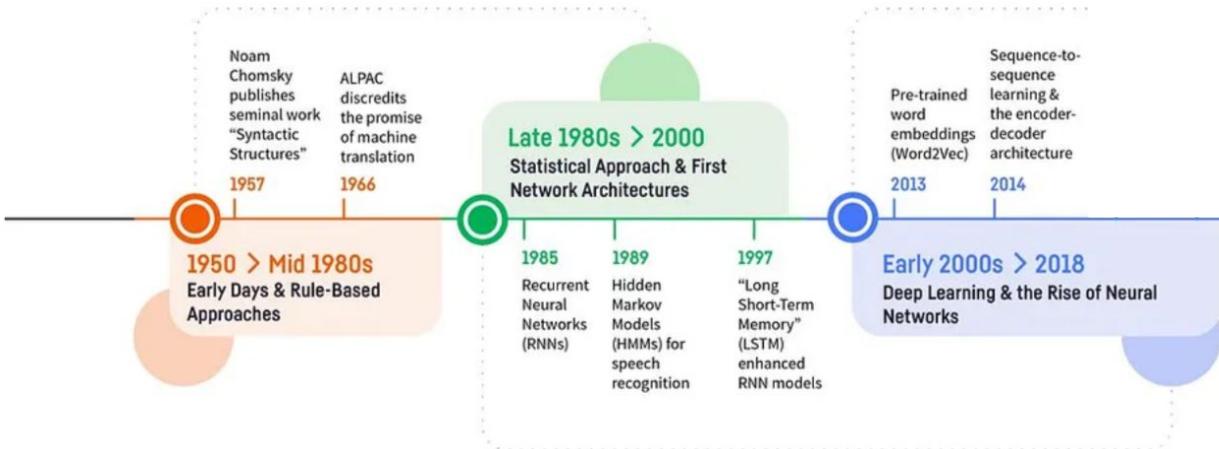


The old times (pre 2017)





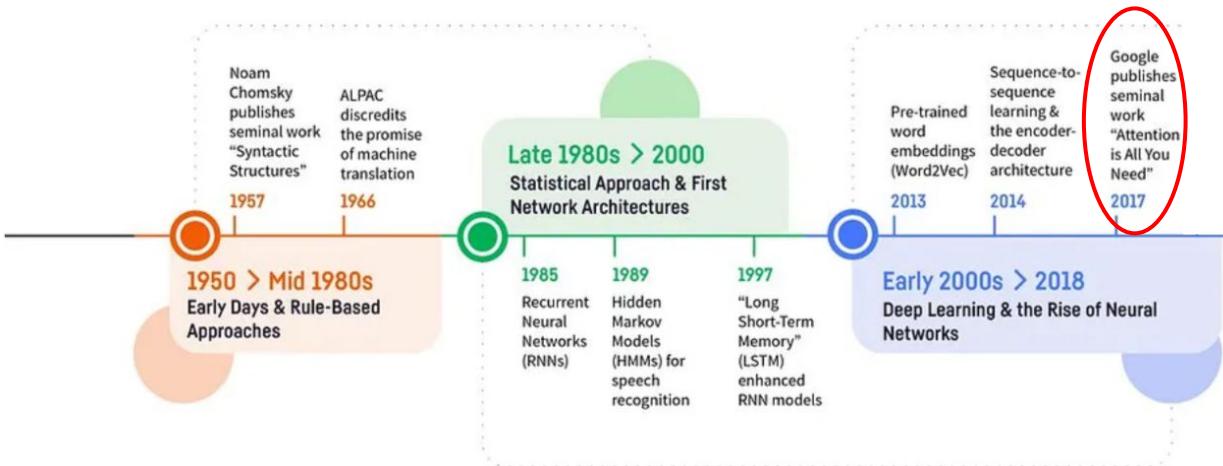
The old times (pre 2017)



- Each problem has its own architecture, with small quirks and problem-specific design choices
- Both in NLP and Computer Vision



The old times (pre 2017)

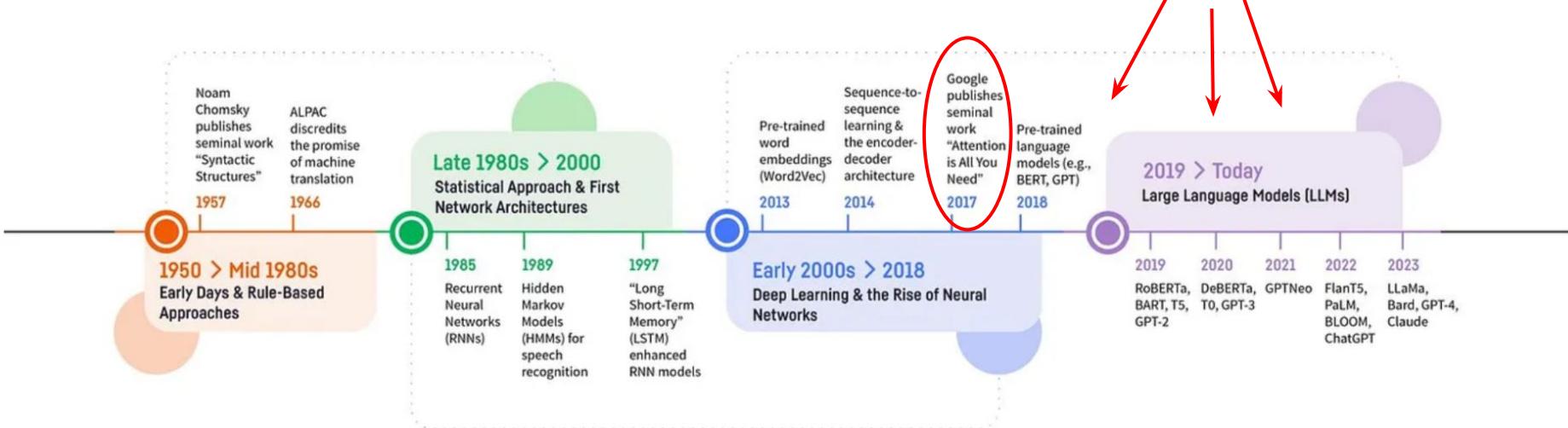


- Each problem has its own architecture, with small quirks and problem-specific design choices
- Both in NLP and Computer Vision
- 2017 - Google published “**Attention is all you need**”



Now (post 2017)

Almost all models are transformers!



- Each problem has its own architecture, with small quirks and problem-specific design choices
- Both in NLP and Computer Vision
- 2017 - Google published “**Attention is all you need**”



The story of deep learning in a sentence

Old methods, developed in the '80s and '90s, when scaled up, started to work.



"We think the most benefits will go to whoever has the biggest computer."



Greg Brockman, OpenAI's CTO



“Compute is getting cheaper faster than we are becoming better researchers”



Hyung Won Chung
Research Scientist at OpenAI



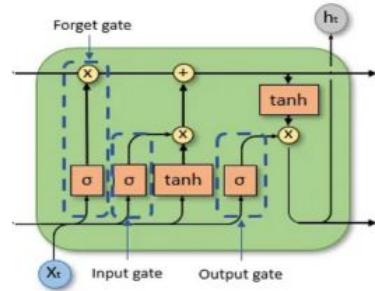
Why now?

Past



+

LSTM



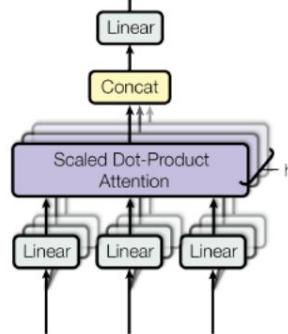
= Doesn't Work

Present



+

Multi-Head Attention



= Works



Seq2Seq, Google, 2014

- Use LSTM as encoder-decoder
 - Input is a sequence → output is a sequence
 - Solve machine translation → solve language modelling
-
- *“The success of our simple LSTM-based approach on MT suggests that it should do well on many other sequence learning problems, provided they have enough training data.”*

Sequence to Sequence Learning with Neural Networks

Ilya Sutskever
Google
ilyasu@google.com

Oriol Vinyals
Google
vinyals@google.com

Quoc V. Le
Google
qvl@google.com

Abstract

Deep Neural Networks (DNNs) are powerful models that have achieved excellent performance on difficult learning tasks. Although DNNs work well whenever large labeled training sets are available, they cannot be used to map sequences to sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a multilayered Long Short-Term Memory (LSTM) to map the input sequence to a vector of a fixed dimensionality, and then another deep LSTM to decode the target sequence from the vector. Our main result is that on an English to French translation task from the WMT’14 dataset, the translations produced by the LSTM achieve a BLEU score of 34.8 on the entire test set, where the LSTM’s BLEU score was penalized on out-of-vocabulary words. Additionally, the LSTM did not have difficulty on long sentences. For comparison, a phrase-based SMT system achieves a BLEU score of 33.3 on the same dataset. When we used the LSTM to rerank the 1000 hypotheses produced by the aforementioned SMT system, its BLEU score increases to 36.5, which is close to the previous best result on this



Pretraining on text - Next Token Prediction

- Given some words, what comes word next?
- Examples:
 - One, two, three, four ... [?]



Pretraining on text - Next Token Prediction

- Given some words, what comes word next?
- Examples:
 - One, two, three, four ... [?]
 - The capital of France is ... [?]



Pretraining on text - Next Token Prediction

- Given some words, what comes word next?
- Examples:
 - One, two, three, four ... [?]
 - The capital of France is ... [?]
 - $25 + 13 = \dots$ [?]



Pretraining on text - Next Token Prediction

- Given some words, what comes word next?
- Examples:
 - One, two, three, four ... [?]
 - The capital of France is ... [?]
 - $25 + 13 = \dots$ [?]
 - $\int \frac{1}{x} dx = \dots$ [?]



Pretraining on text - Next Token Prediction

- Given some words, what comes word next?
- Examples:
 - One, two, three, four ... [?]
 - The capital of France is ... [?]
 - $25 + 13 = \dots$ [?]
 - $\int \frac{1}{x} dx = \dots$ [?]
 - The solution to the Riemann Hypothesis is ... [?]

A form of **inductive reasoning**





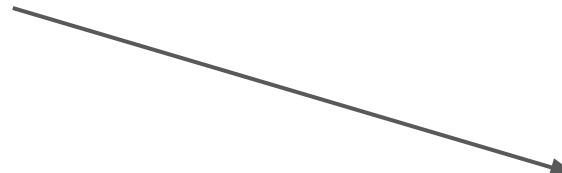
Pretraining on text - Next Token Prediction

- Given some words, what comes word next?
 - Train the model to do this for every word in a sentence.
 - [?] is computed as a classification problem over the vocabulary
-
- The ... [?]
 - The capital ... [?]
 - The capital of ... [?]
 - The capital of France ... [?]
 - The capital of France is ... [?]
- A form of (pretext-based) self-supervised learning!

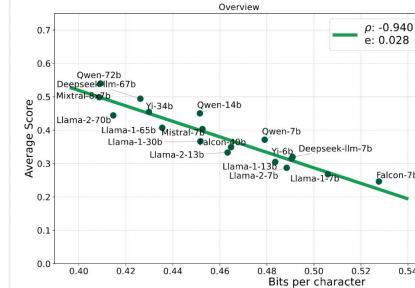


Next Token Prediction

- Simple to describe and implement
- Difficult for the model to do accurately
- For the model to accurately predict what comes next, it needs to “understand” the world
 - I.e., **compress** the information in the corpus



Intelligence = Compression



Compression Represents Intelligence Linearly

Yuchen Huang¹, Jinghan Zhang^{1*}, Zifei Shan², Junstan He²
¹The Hong Kong University of Science and Technology ²Tencent
{yhuang53, jzhang1v, junxianh}@ace.ust.hk

Abstract

There is a belief that learning to compress well will lead to intelligence (Hutter, 2006). Recently, language modeling has been shown to be equivalent to compression, offering a new perspective for the development of large language models (LLMs): the development of more advanced language models is essentially entailing compression which facilitates intelligence. Despite this promising direction, there is little quantitative assessment for the interplay between compression and intelligence. In this work, we examine their relationship in the context of LLMs, treating LLMs as data compression engines. We find that the relationship between the average downstream benchmark scores as a surrogate, specifically targeting intelligence related to knowledge and commonsense, coding, and mathematical reasoning, and the compression efficiency of 20 public LLMs that originate from diverse organizations. Remarkably, we find that LLM intelligence – reflected by average benchmark scores – almost linearly decreases as the bits per character increase. These results provide concrete evidence supporting the belief that superior compression indicates greater intelligence. Furthermore, our findings suggest that our proposed metric, the average downstream benchmark scores on raw text corpora, serves as a reliable evaluation measure that is linearly associated with the model capabilities. We open-source our compression datasets as well as our data collection pipelines to facilitate future researchers to assess compression properly.



2015 - Ideas were already in the collective consciousness

Andrej Karpathy blog

About

The Unreasonable Effectiveness of Recurrent Neural Networks

May 21, 2015

There's something magical about Recurrent Neural Networks (RNNs). I still remember when I trained my first recurrent network for [Image Captioning](#). Within a few dozen minutes of training my first baby model (with rather arbitrarily-chosen hyperparameters) started to generate very nice looking descriptions of images that were on the edge of making sense. Sometimes the ratio of how simple your model is to the quality of the results you get out of it blows past your expectations, and this was one of those times. What made this result so shocking at the time was that the common wisdom was that RNNs were supposed to be difficult to train (with more experience I've in fact reached the opposite conclusion). Fast forward about a year: I'm training RNNs all the time and I've witnessed their power and robustness many times, and yet their magical outputs still find ways of amusing me. This post is about sharing some of that magic with you.

We'll train RNNs to generate text character by character and ponder the question "how is that even possible?"

By the way, together with this post I am also releasing [code on Github](#) that allows you to train character-level language models based on multi-layer LSTMs. You give it a large chunk of text and it will learn to generate text like it one character at a time. You can also use it to reproduce my experiments below. But we're getting ahead of ourselves; What are RNNs anyway?

Recurrent Neural Networks

Sequences. Depending on your background you might be wondering: *What makes Recurrent Networks so special?* A glaring limitation of Vanilla Neural Networks (and also Convolutional Networks) is that their API is too constrained: they accept a fixed-sized vector as input (e.g. an image) and produce a fixed-sized vector as output

Linux Source Code

I wanted to push structured data to its limit, so for the final challenge I decided to use code. In particular, I took all the source and header files found in the [Linux repo on Github](#), concatenated all of them in a single giant file (474MB of C code) (I was originally going to train only on the kernel but that by itself is only ~16MB). Then I trained several as-large-as-my-GPU 3-layer LSTMs over a period of a few days. These models have about 10 million parameters, which is still on the lower end for RNN models. The results are superfun:

```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */

static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            block_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = __in_SB(in.addr);
    selector = seg / 16;
    setup.works = true;
    for (i = 0; i < blocks; i++) {
        seg = i * 16;
        bdp = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
        rws->name = "Getjbreg";
        bprm_self_clear(&s)->version);
        reg->new = blocks|(BPF_STATS << info->historidac)) | PFMR_CLOBATHING_SECONDS << 12;
        return segtable;
    }
}
```

Wikipedia

We saw that the LSTM can learn to spell words and copy general syntactic structures. Let's further increase the difficulty and train on structured markdown. In particular, let's take the [Hutter Prize](#) 100MB dataset of raw Wikipedia and train an LSTM. Following [Graves et al.](#), I used the first 96MB for training, the rest for validation and ran a few models overnight. We can now sample Wikipedia articles! Below are a few fun excerpts. First, some basic markdown output:

Naturalism and decision for the majority of Arab countries' capitalide was grounded by the Irish language by [[John Blair]], [[An Imperial Japanese Revolt]], associated with Guangzham's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Imminences]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25|21]] to note, the Kingdom of Costa Rica, unsuccessful fashioned the [[Thrales]], [[Cynth's Dajoard]], known in western [[Scotland]], near Italy to the conquest of India with the conflict. Copyright was the succession of independence in the stop of Syrian influence that was a famous German movement based on a more popular servitious, non-doctrinal and sexual power post. Many governments recognize the military housing of the [[Civil Liberalization and Infantry Resolution 265 National Party in Hungary]], that is sympathetic to be the [[Punjab Resolution]] (PJS)[<http://www.human.yahoo.com/guardian>]. cfM/4800786d1755196389.htm Official economics Adjoint for the Nazism, Montgomery was swear to advance to the resources for those Socialism's rule, was starting to signing a major tripad of aid exile.]]

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that the plumb and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of retreat. The French army had been marching for three days and demanded--namely, simply to follow the enemy up. The French crowd fled at night, the Russians, who had been marching for three days and were reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made to defend itself as by the fact that the Russian bridges broke down, unarmed soldiers, people from Moscow and women with children who had been captured and who had been taken to the ice-covered water and did not, pressed forward into boats and into the ice-covered water and did not,

Get the turns on inside quotes:

You meant to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give the French a hard time. I mean, I can supply you with the finest word he spoke to prove his own rectitude and therefore caused Kutuzov to be annihilated by his own desire.

KUTUZOV: I mean exactly today what I said.

Cell that robustly activates inside if statements:

```
static int dequeue_signals(struct sigpending *pending, sigset_t *mask,
                           siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig >= 0) {
        if (current->notifier) {
            if ((current->notifier->mask & sig) && !((current->notifier->current->notifier_data))) {
                clear_thread_flag(TIE_SIGPENDING);
                return sig;
            }
        }
    }
    collect_signals(sig, pending, info);
    return sig;
}
```



The sentiment neuron - OpenAI, 2017

- Train byte-level LSTM to generate IMDB reviews, unsupervised
 - 4 Pascal GPUs for a month
-
- “we find a single unit which performs sentiment analysis”
 - A particular neuron activates when the review is positive and deactivates when the review is negative

Learning to Generate Reviews and Discovering Sentiment

Alec Radford¹ Rafal Jozefowicz¹ Ilya Sutskever¹

Abstract

We explore the properties of byte-level recurrent language models. When given sufficient amounts of capacity, training data, and compute time, the representations learned by these models include disentangled features corresponding to high-level concepts. Specifically, we find a single unit which performs sentiment analysis. These representations, learned in an unsupervised manner, achieve state of the art on the binary subset of the Stanford Sentiment Treebank. They are also very data efficient. When using only a handful of labeled examples, our approach matches the performance of strong baselines trained on full datasets. We also demonstrate the sentiment unit has a direct influence on the generative process of the model. Simply fixing its value to be positive or negative generates samples with the corresponding positive or negative sentiment.

it is now commonplace to reuse these representations on a broad suite of related tasks - one of the most successful examples of transfer learning to date ([Quab et al., 2014](#)).

There is also a long history of unsupervised representation learning ([Olshausen & Field, 1997](#)). Much of the early research into modern deep learning was developed and validated via this approach ([Hinton & Salakhutdinov, 2006](#)) ([Huang et al., 2007](#)) ([Vincent et al., 2008](#)) ([Coates et al., 2010](#)) ([Le, 2013](#)). Unsupervised learning is promising due to its ability to scale beyond only the subsets and domains of data that can be cleaned and labeled given resource, privacy, or other constraints. This advantage is also its difficulty. While supervised approaches have clear objectives that can be directly optimized, unsupervised approaches rely on proxy tasks such as reconstruction, density estimation, or generation, which do not directly encourage useful representations for specific tasks. As a result, much work has gone into designing objectives, priors, and architectures meant to encourage the learning of useful representations. We refer readers to [Goodfellow et al. \(2016\)](#) for a detailed



The sentiment neuron - OpenAI, 2017

- Train byte-level LSTM to generate IMDB reviews, unsupervised
- 4 Pascal GPUs for a month
- “we find a single unit which performs sentiment analysis”
 - A particular neuron activates when the review is positive and deactivates when the review is negative

Judy Holliday struck gold in 1950 with George Cukor's film version of "Born Yesterday," and from that point forward, her career consisted of trying to find material good enough to allow her to strike gold again. It never happened. In "It Should Happen to You" (I can't think of a blander title, by the way), Holliday does yet one more variation on the dumb blonde who's maybe not so dumb after all, but everything about this movie feels warmed over and half hearted. Even Jack Lemmon, in what I believe was his first film role, can't muster up enough energy to enliven this recycled comedy. The audience knows how the movie will end virtually from the beginning, so mostly it just sits around waiting for the film to catch up. Maybe if you're enamored of Holliday you'll enjoy this; otherwise I wouldn't bother. Grade: C

Once in a while you get amazed over how BAD a film can be, and how in the world anybody could raise money to make this kind of crap. There is absolutely NO talent included in this film - from a crappy script, to a crappy story to crappy acting. Amazing...

Team Spirit is maybe made by the best intentions, but it misses the warmth of "All Stars" (1997) by Jean van de Velde. Most scenes are identical, just not that funny and not that well done. The actors repeat the same lines as in "All Stars" but without much feeling.

God bless Randy Quaid...his lecherous Cousin Eddie in Vacation and Christmas Vacation hilariously stole the show. He even made the awful Vegas Vacation at least worth a look. I will say that he tries hard in this made for TV sequel, but that the script is so NON funny that the movie never really gets anywhere. Quaid and the rest of the returning Vacation vets (including the original Audrey, Dana Barron) are wasted here. Even European Vacation's Eric Idle cannot save the show in a brief cameo.... Pathetic and sad...actually painful to watch....Christmas Vacation 2 is the worst of the vacation franchise.



The sentiment neuron - OpenAI, 2017

*“... our work encourages further research into **language modelling** as it demonstrates that the standard language modelling objective
with no modifications*

is sufficient to learn high-quality representations ...”



Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* †

University of Toronto

aidan@cs.toronto.edu

Łukasz Kaiser*

Google Brain

lukaszkaiser@google.com

Illia Polosukhin* ‡

illia.polosukhin@gmail.com



Attention Is All You Need

Adept

Ashish Vaswani*
Google Brain
avaswani@google.com

character.ai

Noam Shazeer*
Google Brain
noam@google.com

Adept

Niki Parmar*
Google Research
nikip@google.com

 Inceptive

Jakob Uszkoreit*
Google Research
usz@google.com



Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

co:here

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

NEAR
INCORPORATED





Transformer Decoder for Next Token Prediction

Input

<start>

Decoder

Transformer Decoder



Transformer Decoder for Next Token Prediction

Input

<start>

<pred>

Decoder

Transformer Decoder



Transformer Decoder for Next Token Prediction

Input

<start>

Target

<pred>

The

Decoder

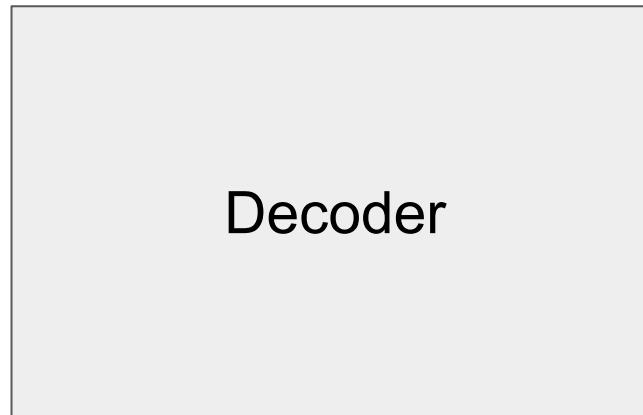
Transformer Decoder



Transformer Decoder for Next Token Prediction

Input

<start>



Target

<pred>

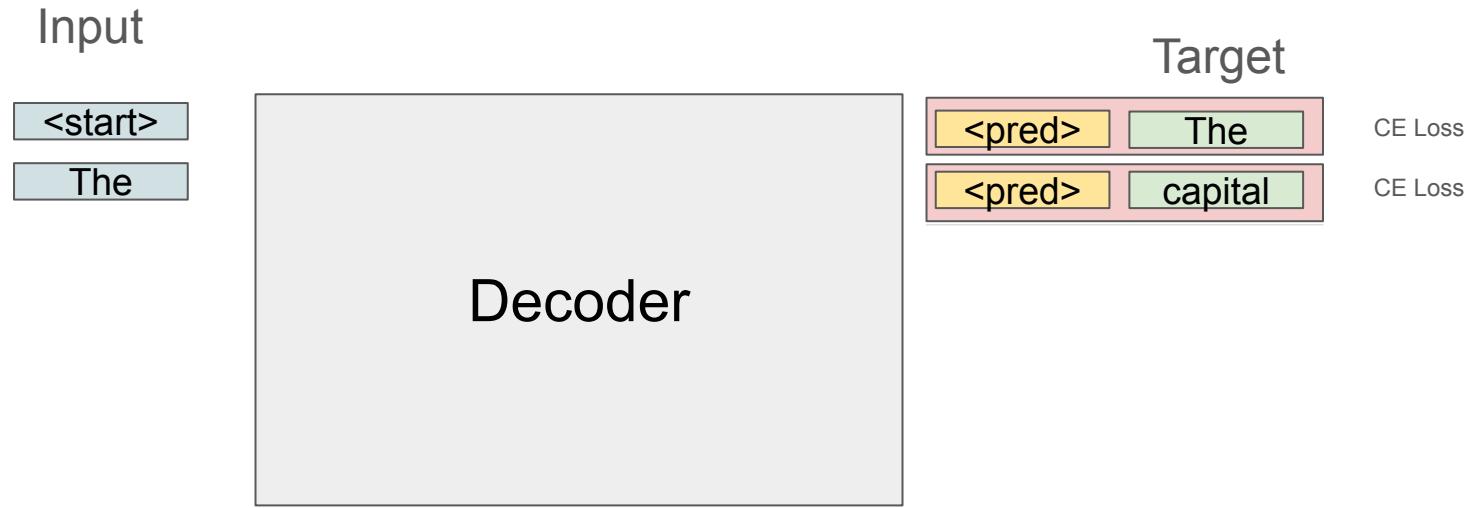
The

CE Loss

Transformer Decoder



Transformer Decoder for Next Token Prediction



Transformer Decoder



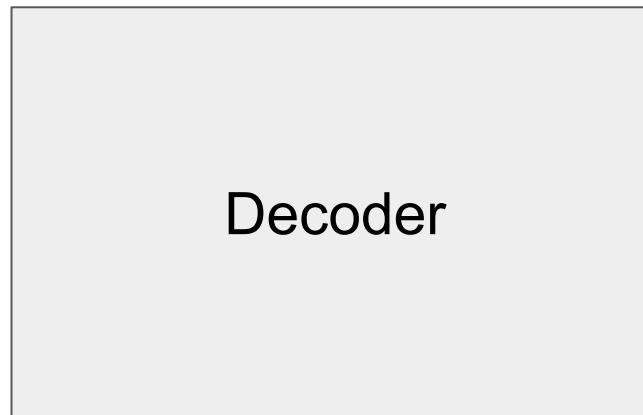
Transformer Decoder for Next Token Prediction

Input

<start>

The

capital



Target

<pred>

The

<pred>

capital

<pred>

of

CE Loss

CE Loss

CE Loss

Transformer Decoder



Transformer Decoder for Next Token Prediction

Input

<start>

The

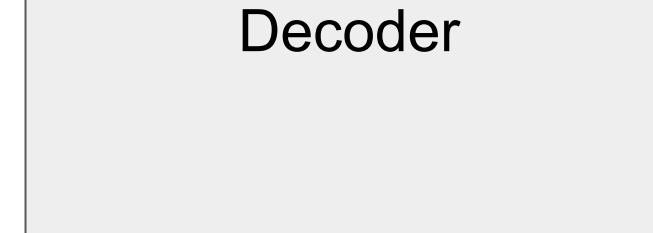
capital

of

France

is

Paris



Target

<pred>

The

CE Loss

<pred>

capital

CE Loss

<pred>

of

CE Loss

<pred>

France

CE Loss

<pred>

is

CE Loss

<pred>

Paris

CE Loss

<pred>

<end>

CE Loss

Transformer Decoder



Transformer Decoder for Next Token Prediction

Input

<start>

The

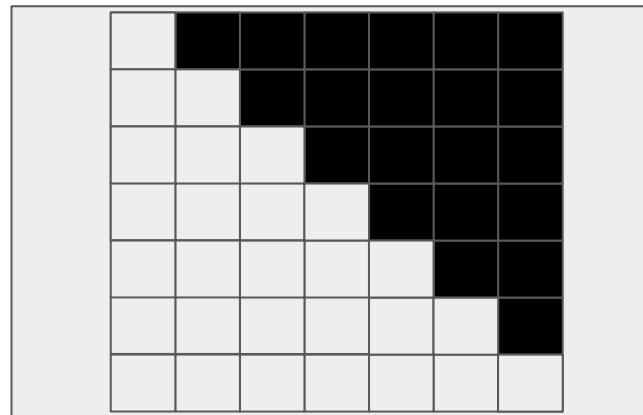
capital

of

France

is

Paris



Target

<pred>

The

CE Loss

<pred>

capital

CE Loss

<pred>

of

CE Loss

<pred>

France

CE Loss

<pred>

is

CE Loss

<pred>

Paris

CE Loss

<pred>

<end>

CE Loss

Using a causal attention mask, we can do this in a single forward pass!

Transformer Decoder



Seq2Seq is all you need?

- People started to realize that all NLP problems can be cast as text in → text out
- Example: Text-to-Text Transfer Transformer (T5), 2019

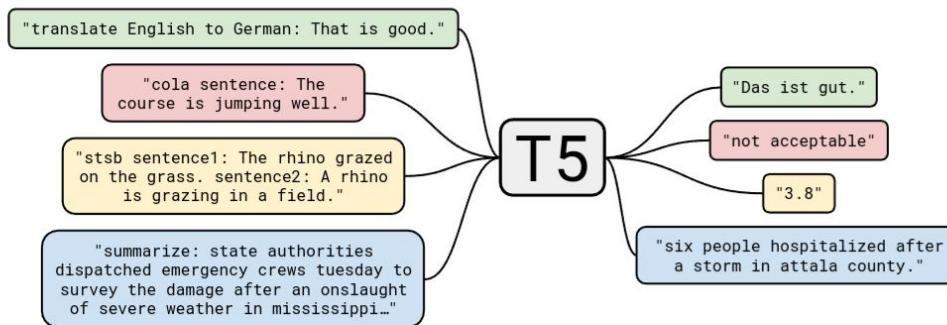


Figure 1: A diagram of our text-to-text framework. Every task we consider—including translation, question answering, and classification—is cast as feeding our model text as input and training it to generate some target text. This allows us to use the same model, loss function, hyperparameters, etc. across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



An aside: Tokenization

- Text tokenization = splitting a string into a sequence of tokens
- Most methods so far used either
 - Word-level tokenization
 - Character-level tokenization
- Both are inadequate
 - Splitting by words results in many “out-of-vocabulary” words
 - Misspellings? Many edge cases.
 - Splitting text by characters results in very large sequence length which is computationally intractable



Tokenization: Word-level splitting

- What happens when we encounter a word that we have never seen in our training data?
 - Not much we can do
 - Assign a special <UNK> token
 - Lose a lot of meaning
 - Especially hurts in texts/languages with many rare words/entities

- “amazing!”, “state-of-the-art”, “un-thinkable”, “prize-winning”, “aren’t”, “O’Neill”
- Some languages don’t even use spaces to mark word boundaries!



Tokenization: Character-level splitting

- Vocabulary size is small
 - 256 entries
 - No unknown tokens
-
- However, sequence is much larger
 - Need to learn from scratch how to combine characters into words



Byte-Pair-Encoding Tokenization: The middle ground

- Tokenize into subwords using BPE
 - A **greedy** compression algorithm
 - Vocabulary size fixed and specified by us
 - No unknown words
 - Always fallback to small subwords
- $\mathcal{V} \leftarrow$ All characters in the training data (as base **tokens**)
 - For k steps:
 - Tokenize the data, taking the longest prefix each time
 - Count the frequency of adjacent token pairs in the data
 - Choose the pair $\langle l, r \rangle$ that occurs most frequently
 - Add the pair to the vocabulary as a new token $\mathcal{V} \leftarrow \mathcal{V} \cup \{lr\}$
 - Return \mathcal{V}



BPE - Example

GPT-4o & GPT-4o mini

GPT-3.5 & GPT-4

GPT-3 (Legacy)

Ilya Sutskever FRS (born 8 December 1986) is a Canadian-Israeli-Russian computer scientist who specializes in machine learning

Clear

Show example

Tokens

29

Characters

126

Ilya Sutskever FRS (born 8 December 1986) is a Canadian-Israeli-Russian computer scientist who specializes in machine learning

[40, 306, 64, 311, 5500, 365, 332, 376, 6998, 357, 6286, 807, 3426, 12113, 8, 318, 257, 5398, 12, 29818, 12, 16220, 3644, 11444, 508, 29786, 287, 4572, 4673]



HuggingFace tokenizers in Rust



- Official tokenizers library is built in Rust
- Python bindings
- Anyone who ever used a HuggingFace LLM has used this implementation

```
use tokenizers::tokenizer::{Result, Tokenizer};

fn main() -> Result<()> {
    // needs http feature enabled
    let tokenizer = Tokenizer::from_pretrained("bert-base-cased", None)?;

    let encoding = tokenizer.encode("Hey there!", false)?;
    println!("{:?}", encoding.get_tokens());
    Ok(())
}
```



GPT-1, OpenAI, 2018

- Pretrain decoder-only transformer on BPE tokenized text
 - 12 layers, dmodel = 768
 - Books corpus
 - 12 GPUs for a month
- Pretraining for next-token prediction helps a lot!
- Fine-tune for classification last embedding

Improving Language Understanding by Generative Pre-Training

Alec Radford

OpenAI

alec@openai.com

Karthik Narasimhan

OpenAI

karthikn@openai.com

Tim Salimans

OpenAI

tim@openai.com

Ilya Sutskever

OpenAI

ilyasu@openai.com

Abstract

Natural language understanding comprises a wide range of diverse tasks such as textual entailment, question answering, semantic similarity assessment, and document classification. Although large unlabeled text corpora are abundant, labeled data for learning these specific tasks is scarce, making it challenging for discriminatively trained models to perform adequately. We demonstrate that large gains on these tasks can be realized by *generative pre-training* of a language model on a diverse corpus of unlabeled text, followed by *discriminative fine-tuning* on each specific task. In contrast to previous approaches, we make use of task-aware input transformations during fine-tuning to achieve effective transfer while requiring minimal changes to the model architecture. We demonstrate the effectiveness of our approach on a wide range of benchmarks for natural language understanding. Our general task-agnostic model outperforms discriminatively trained models that use architectures specifically crafted for each task, significantly improving upon the state of the art in 9 out of the 12 tasks studied. For instance, we achieve absolute improvements of 8.9% on commonsense reasoning (Stories Cloze Test), 5.7% on question answering (RACE), and 1.5% on textual entailment (MultiNLI).



February 14, 2019

Better language models and their implications

[Read paper ↗](#)

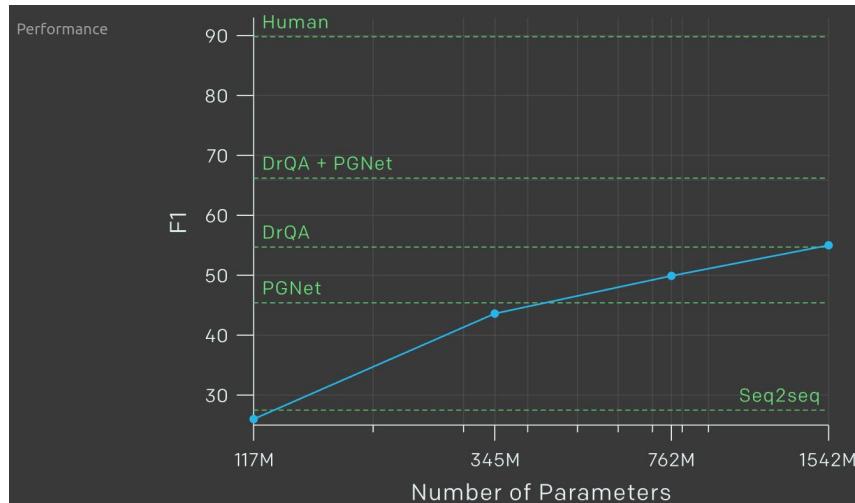
[View code ↗](#)

ur highest priority, ington and Ashland, Hicks said. She was carrying a pair of **Hicks** and **Popple** and **Popple** was taking far too m
t land and **Race**. These four-hour em to be quite comfornable. A
'rombie and Fitch on Holly Accordin to Peter Reissosp'ed that the creatures 'Lines' on the front and 'Fas
hicks also spoke some f
Midwest fashion painted "Warfa
ny people think of the Warfa
research Triangle Park nuclear
ked blue, with some crystals
and 'Fashion Police' on the back.
lack jacket, black jeans
enomenon is finally
onsequences on public an
lack jacket, black jeans
ence who had some really cra
working with the Federal Ra
n of evolution, or at le
where the animals wer
'Our top priority is to secure
some crystals on top,' said R
king with the Federal Railroa
k jacket, black jeans, and b
ur highest priority, Hicks s
mething like a diale
'It's happened the Civil War
d Lines' on the front and 'Fas
hicks with some crystals or
my people in the front and
ation the front and 'Fashion
unded by two peaks of rock,



GPT-2, OpenAI, 2019

- Direct scale-up of GPT-1
- 1.5B parameters
- 40GB of text
 - Web pages that were linked on a reddit posts with at least 3 likes
- More model parameters improves performance





When people started to notice: GPT-3, OpenAI, 2020

- 175B parameters
- 300B tokens
- New capabilities emerged

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	3.0×10^{-4}
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M	0.6×10^{-4}

Table 2.1: Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

Language Models are Few-Shot Learners

Tom B. Brown*	Benjamin Mann*	Nick Ryder*	Melanie Subbiah*	
Jared Kaplan†	Prafulla Dhariwal	Arvind Neelakantan	Pranav Shyam	Girish Sastry
Amanda Askell	Sandhini Agarwal	Ariel Herbert-Voss	Gretchen Krueger	Tom Henighan
Rewon Child	Aditya Ramesh	Daniel M. Ziegler	Jeffrey Wu	Clemens Winter
Christopher Hesse	Mark Chen	Eric Sigler	Mateusz Litwin	Scott Gray
Benjamin Chess	Jack Clark	Christopher Berner		
Sam McCandlish	Alec Radford	Ilya Sutskever	Dario Amodei	

OpenAI

Abstract

Recent work has demonstrated substantial gains on many NLP tasks and benchmarks by pre-training on a large corpus of text followed by fine-tuning on a specific task. While typically task-agnostic in architecture, this method still requires task-specific fine-tuning datasets of thousands or tens of thousands of examples. By contrast, humans can generally perform a new language task from only a few examples or from simple instructions – something which current NLP systems still largely struggle to do. Here we show that scaling up language models greatly improves task-agnostic,



Scaling law: performance follows a power-trend with compute

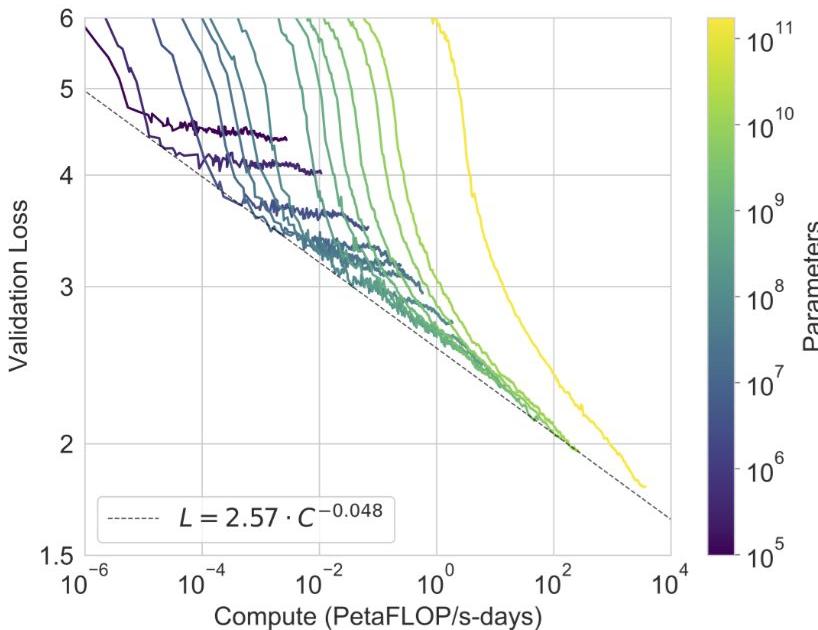
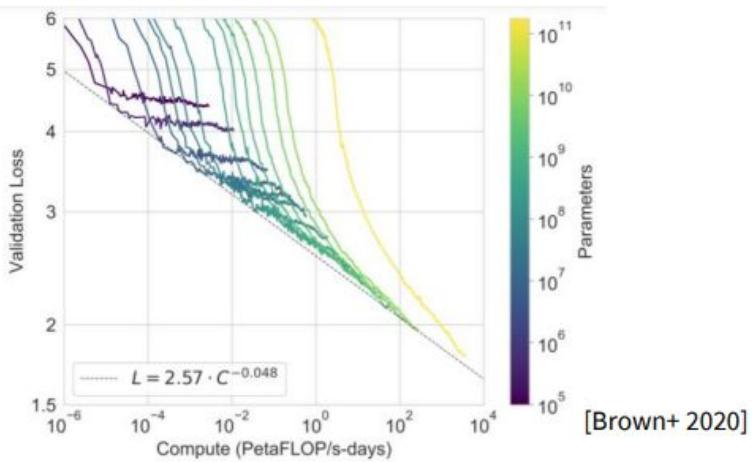
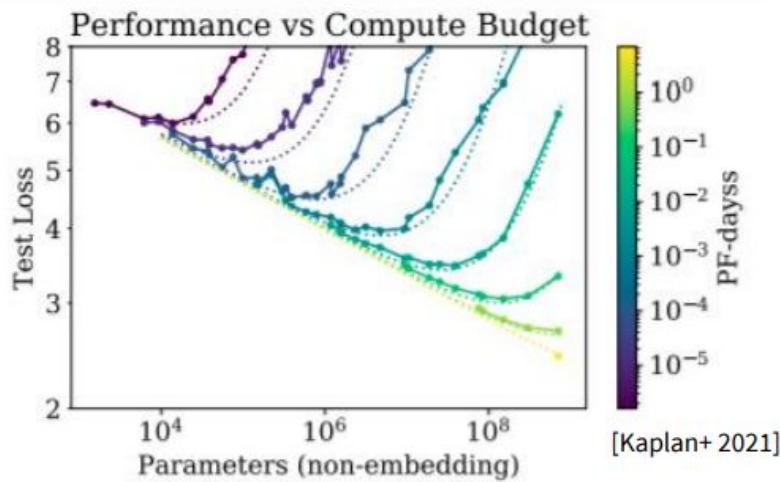


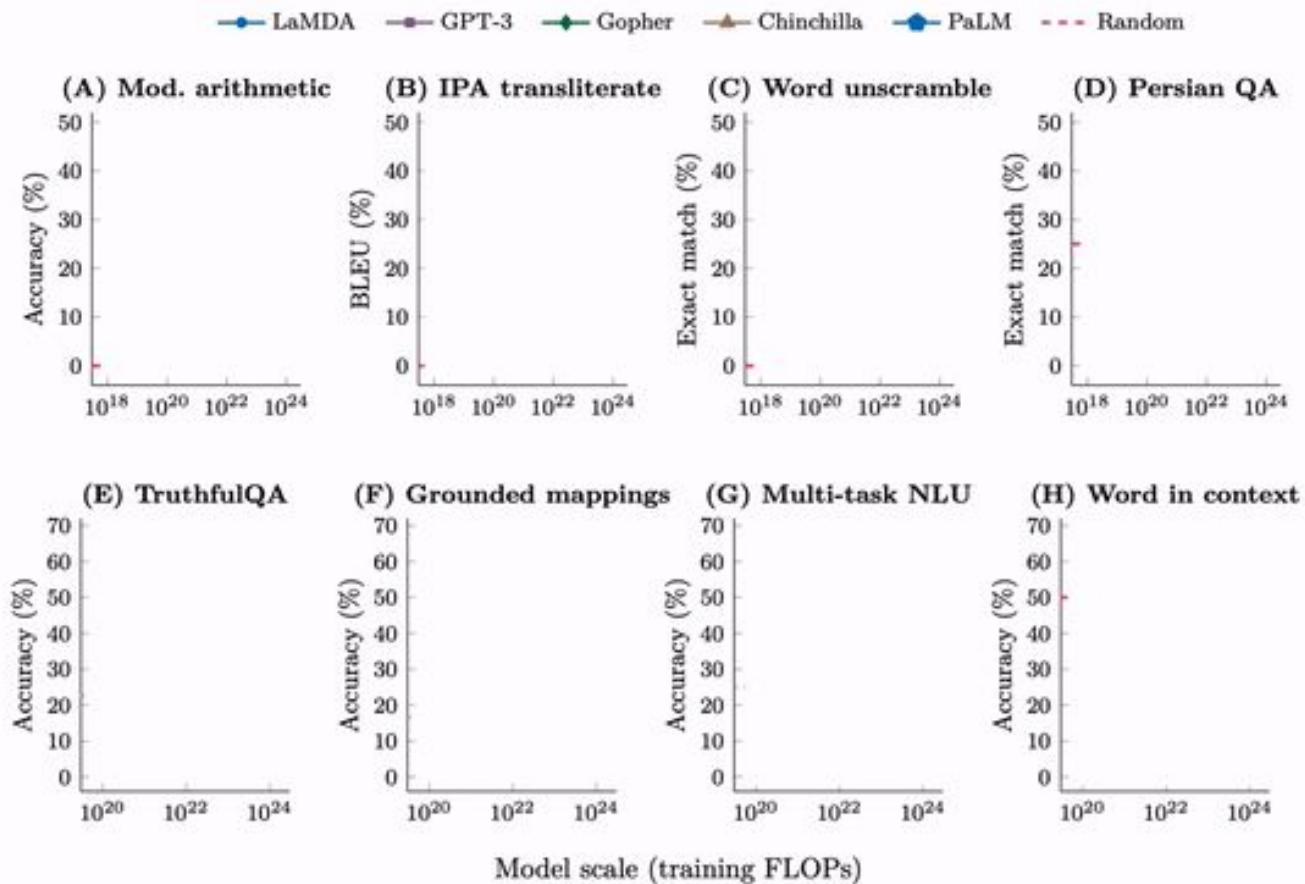
Figure 3.1: Smooth scaling of performance with compute. Performance (measured in terms of cross-entropy validation loss) follows a power-law trend with the amount of compute used for training. The power-law behavior observed in [KMH⁺20] continues for an additional two orders of magnitude with only small deviations from the predicted curve. For this figure, we exclude embedding parameters from compute and parameter counts.



Scaling is all you need!



Bigger is not only better!



Emergent Abilities of Large Language Models

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph et al. (2022)

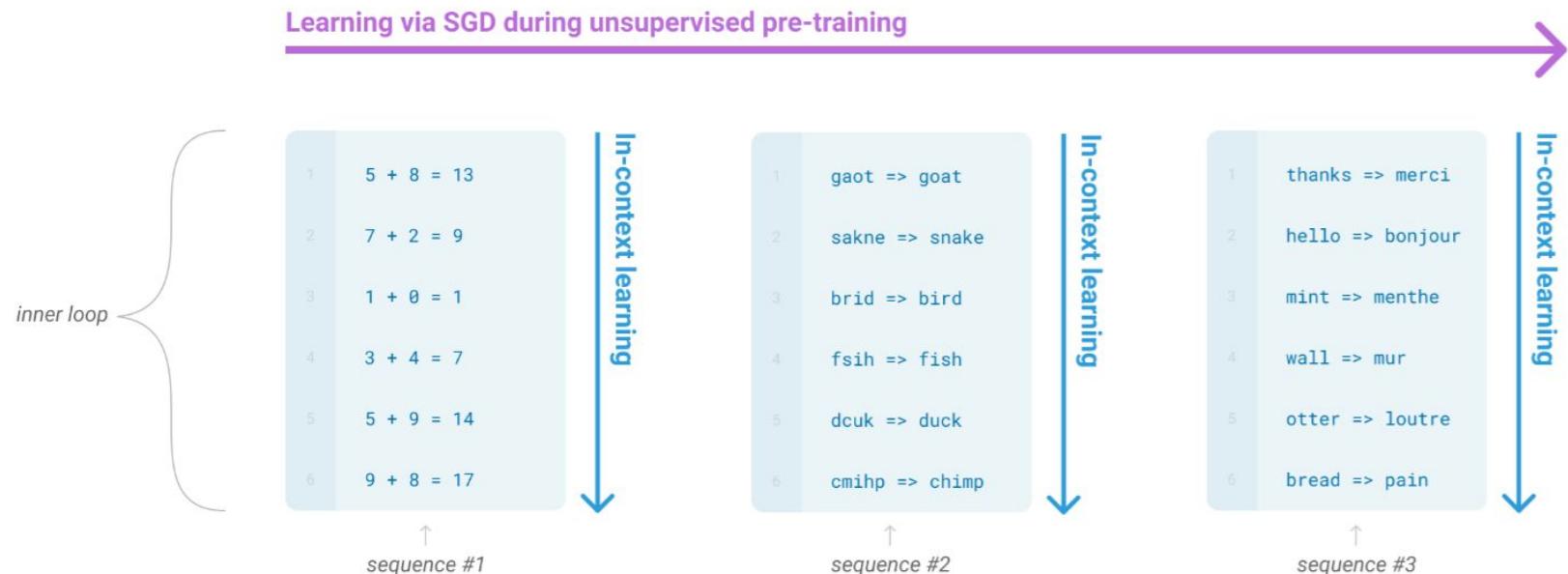


Paradigm shift: In-Context Learning

- Don't need to fine-tune directly, put examples directly in the prompt.
- Paradigm shift.
 - Initially, people got data and trained a model on it
 - Then, general pretrained models (i.e., BERT) were made available, but they were useless
 - People still needed to fine-tune.
 - Now, fine-tuning became “obsolete”,
 - use model as-is, put examples in the prompt



Examples of “in-context-learning” in training data



Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan et al. "Language models are few-shot learners." *Advances in neural information processing systems* 33 (2020): 1877-1901.



“Base-models” don’t follow instructions

Prompt Explain the moon landing to a 6 year old in a few sentences.

Completion GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.



InstructGPT, OpenAI, 2022

January 27, 2022

Aligning language models to follow instructions

[Read paper ↗](#)

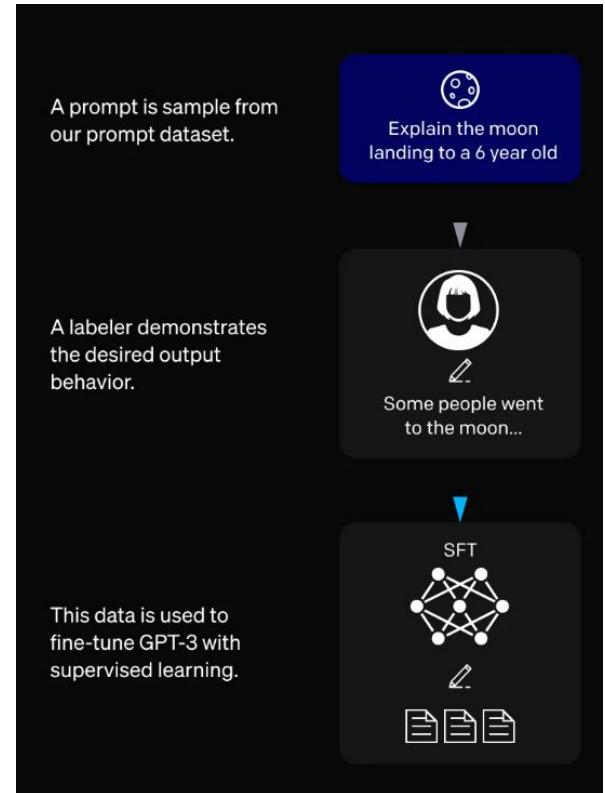
[View model card ↗](#)





Supervised Fine-Tuning (SFT)

- Instruction-following
- Gather a dataset of instruction and human-generated responses
- Train the model to output the response





<start>

What

is

the

capital

of

France?

The

capital

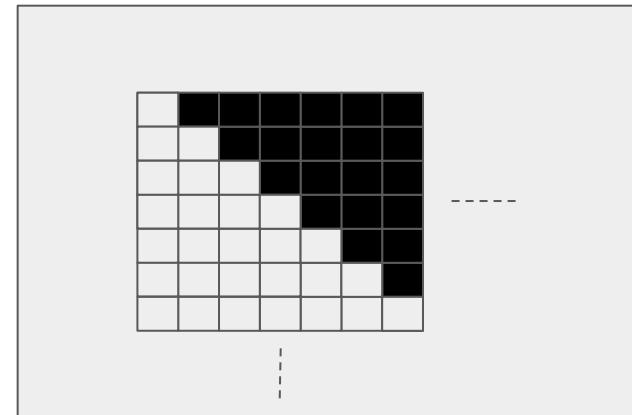
of

France

is

Paris

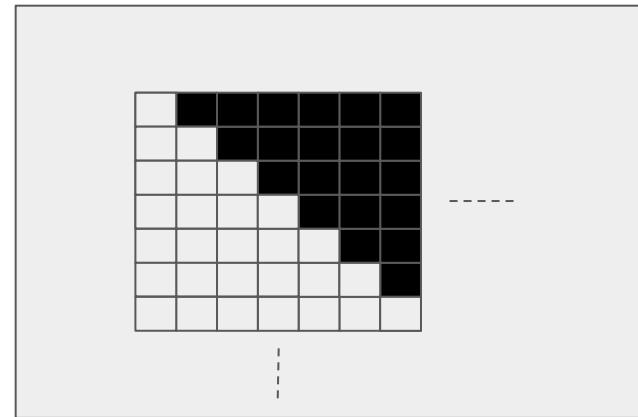
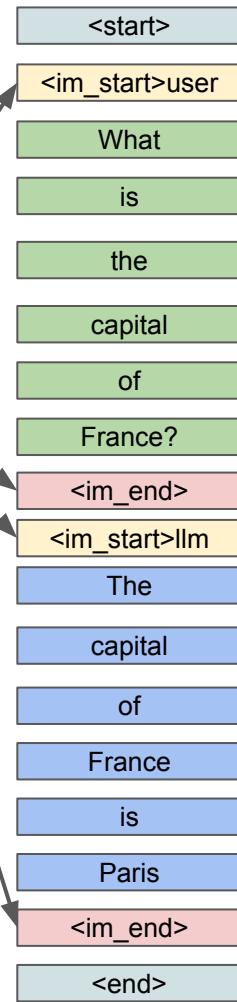
<end>



Base-model
(pretrained transformer decoder)



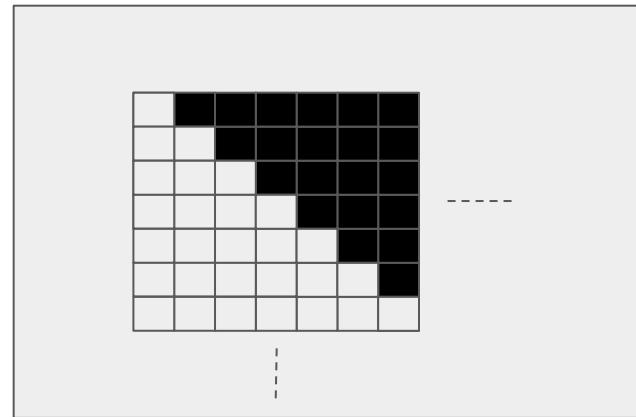
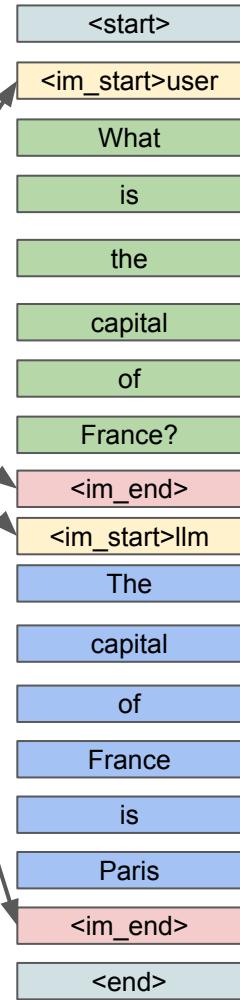
Add
special
tokens



Base-model
(pretrained transformer decoder)



Add
special
tokens



Base-model
(pretrained transformer decoder)



Same as before, next token prediction!



Now, the model can follow instruction (i.e. chat)

Prompt Explain the moon landing to a 6 year old in a few sentences.

Completion GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

InstructGPT

People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

November 30, 2022

Introducing ChatGPT

[Try ChatGPT ↗](#)

[Download ChatGPT desktop >](#)

[Learn about ChatGPT >](#)

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

ChatGPT is a sibling model to [InstructGPT](#), which is trained to follow an instruction in a prompt and provide a detailed response.

We are excited to introduce ChatGPT to get users' feedback and learn about its strengths and weaknesses. During the research preview, usage of ChatGPT is free. Try it now at [chatgpt.com](#).

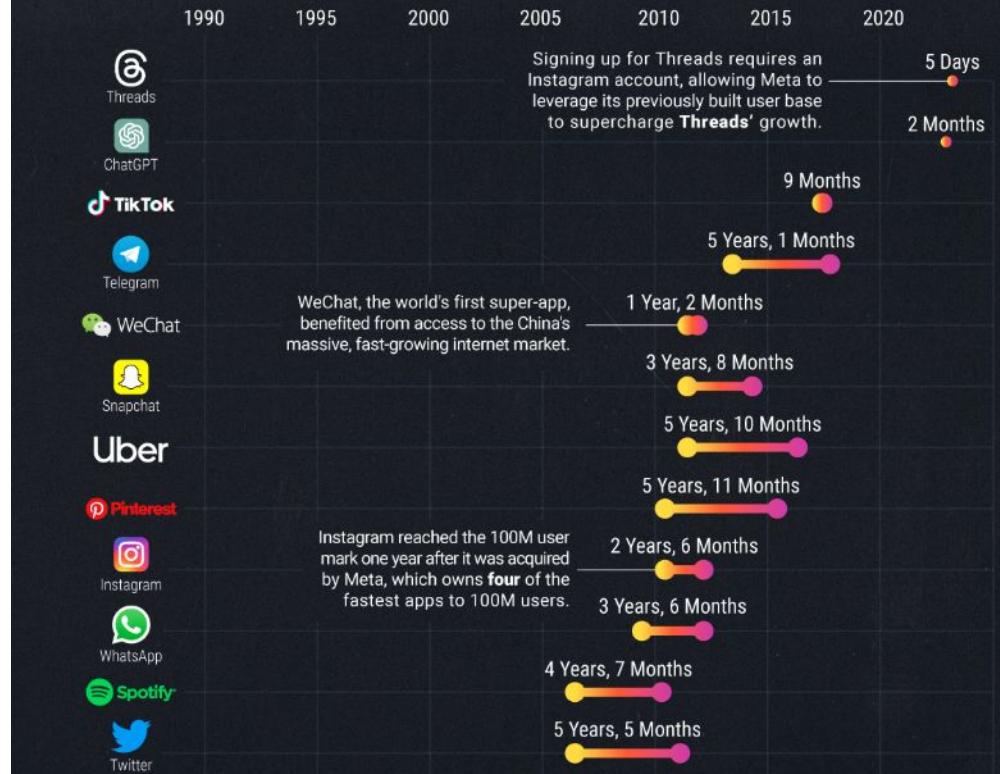
ChatGPT, 2022

- 2 months to reach 100M users
- First time a chatbot can answer questions, refuse to answer, argue with you, “understands” you

HOW LONG DID IT TAKE APPS TO REACH

100M Users?

Meta's newest social media platform, Threads, took less than a week to attract 100 million users to its platform, smashing the previous record of 2 months held by OpenAI's ChatGPT.





Capability prediction on 23 coding problems

– Mean Log Pass Rate

5

4

3

2

1

0

1μ

10μ

100μ

0.001

0.01

0.1

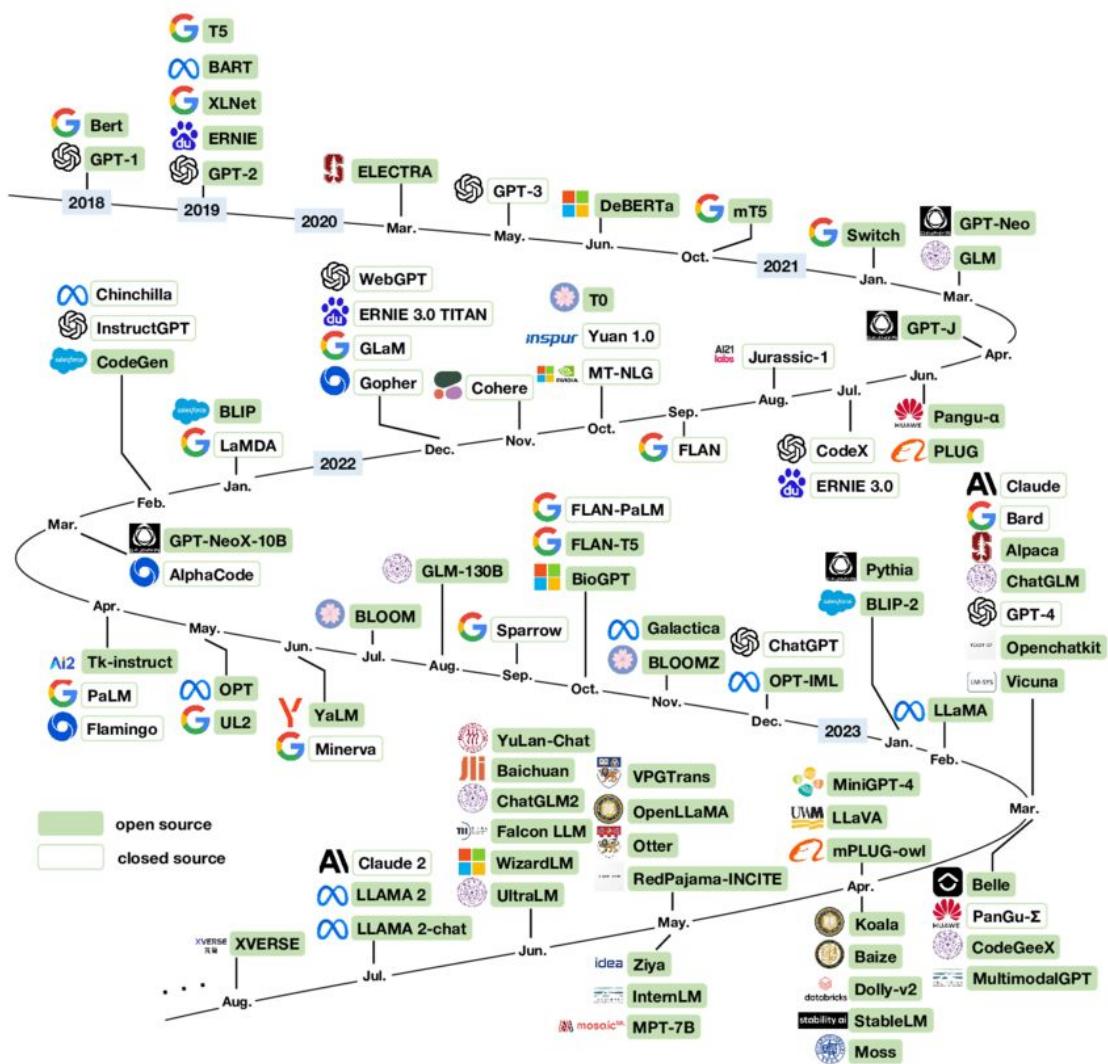
1

Compute

- Observed
- - - Prediction
- gpt-4

Figure 2. Performance of GPT-4 and smaller models. The metric is mean log pass rate on a subset of the HumanEval dataset. A power law fit to the smaller models (excluding GPT-4) is shown as the dotted line; this fit accurately predicts GPT-4's performance. The x-axis is training compute normalized so that GPT-4 is 1.

Cambrian Explosion of LLMs





Ilya Sutskever - NeurIPS talk 2024

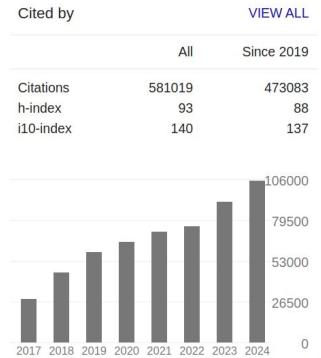
Pre-training as we know it will end

Compute is growing:

- Better hardware
- Better algorithms
- Larger clusters

Data is not growing:

- We have but one internet
- **The fossil fuel of AI**

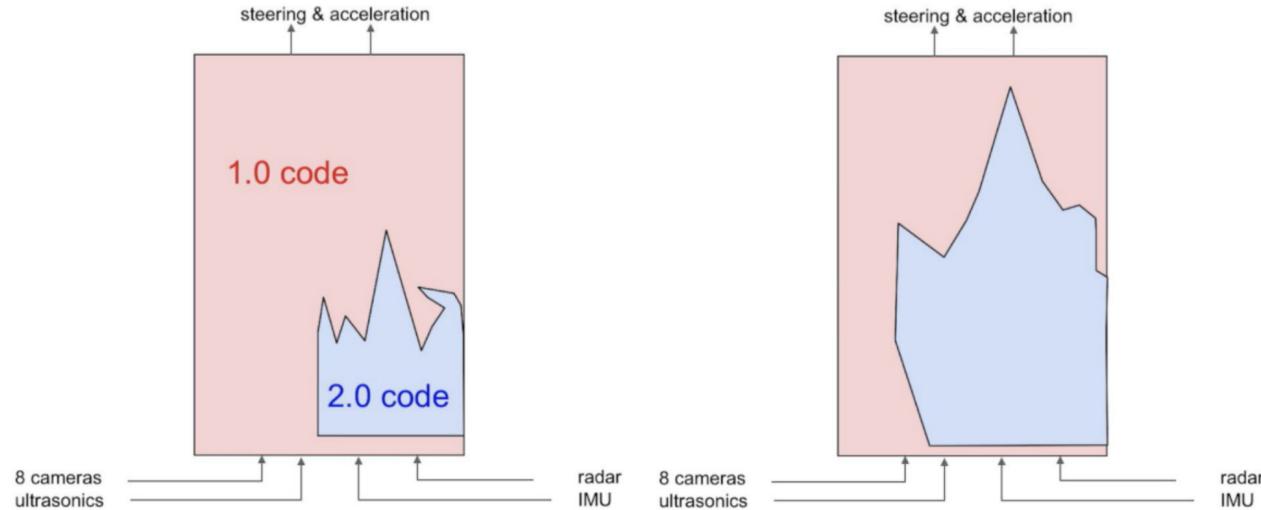




Recall from last lecture ...

Software 1.0: “classical” software engineering

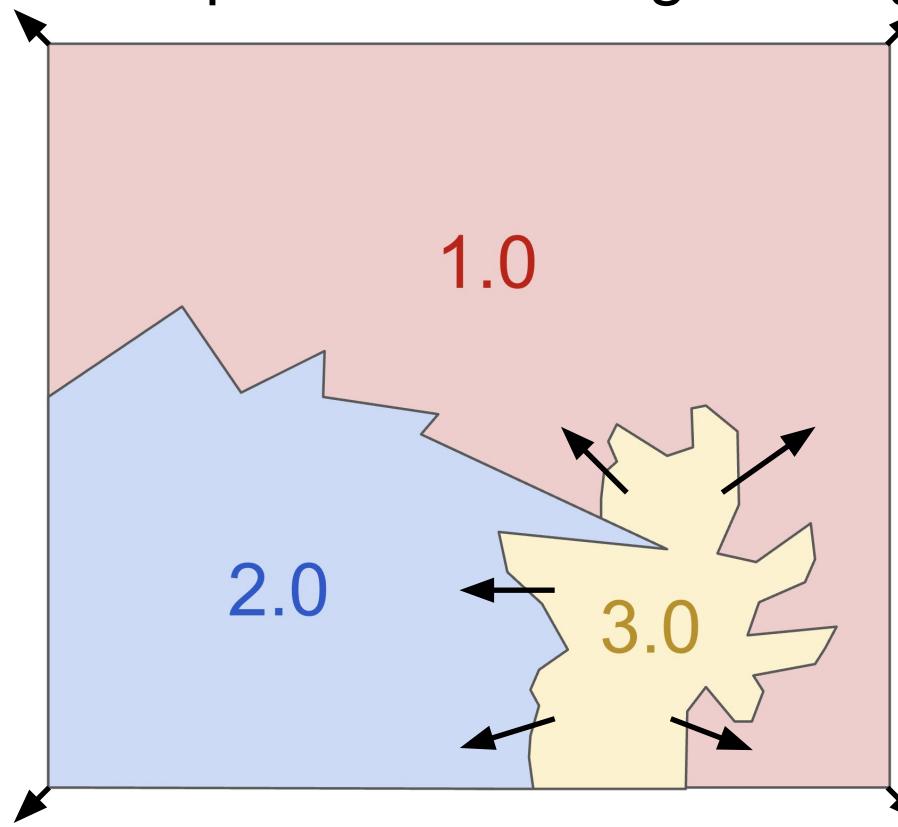
Software 2.0: (neural network) weights



Source: Andrej Karpathy



Software 3.0: Prompt / Context Engineering



Source: Andrej Karpathy

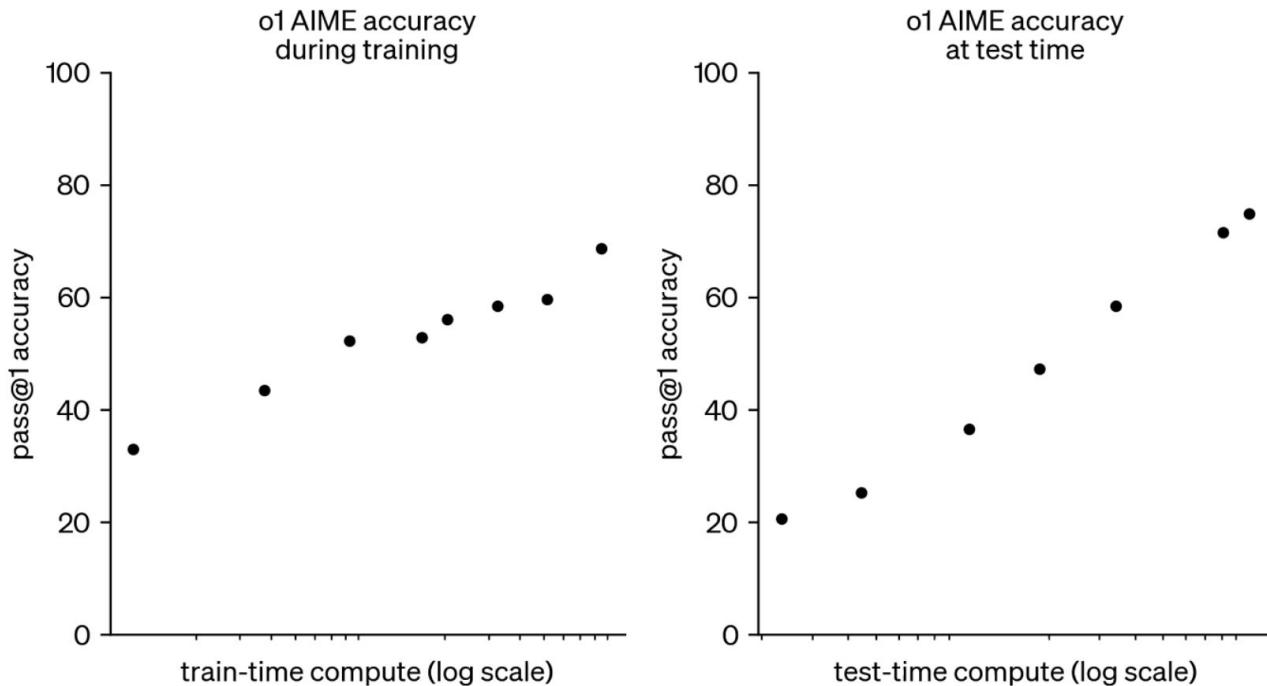


Test-time compute

- We can trade test-time compute for model size
- Bigger models can store more knowledge
- But do we really need to store knowledge in the model? What if we can get answers by thinking / reasoning more vs memorizing?
- What if we can retrieve knowledge from an external source on demand?
- For on-device LLMs, ideally what we want is a “**Small Reasoning Model**” (SLM)



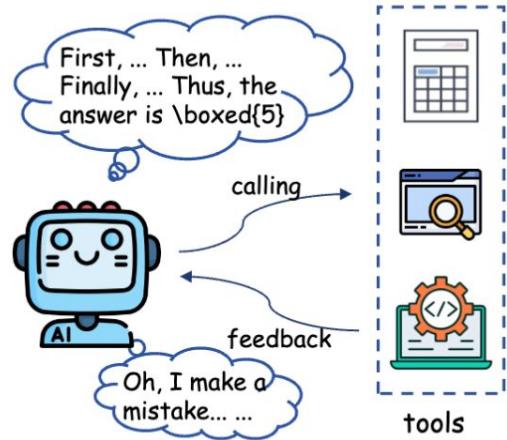
We can “trade” training-time compute with test-time compute



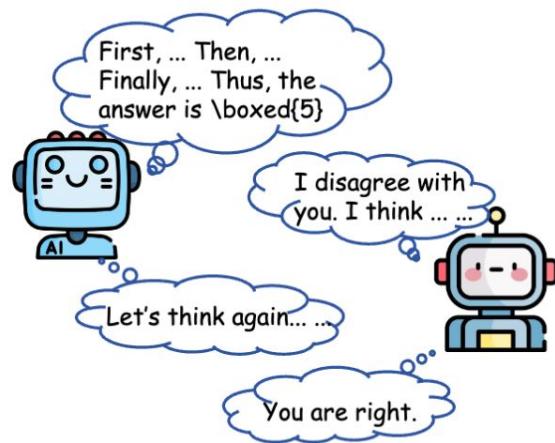
The impact of compute time on OpenAI’s o1 model accuracy during training (left) and test time (right), highlighting the effectiveness of Test Time Compute. Credit: OpenAI.



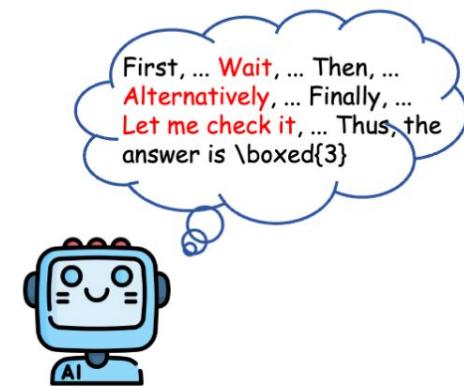
Other types of test-time compute



(a) Tool checking



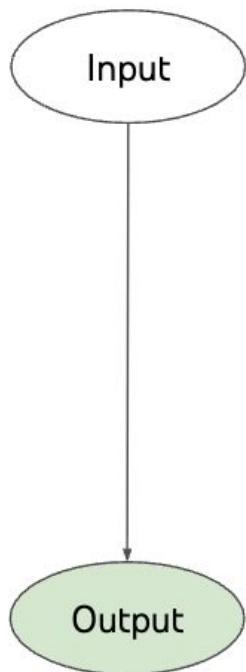
(b) External model evaluation



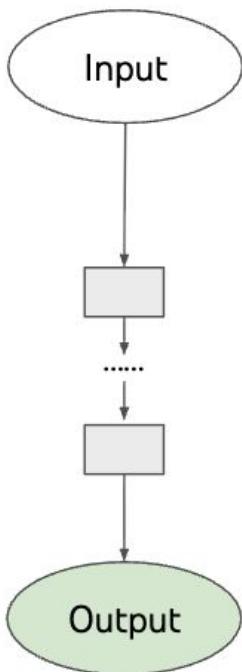
(c) Self-critique



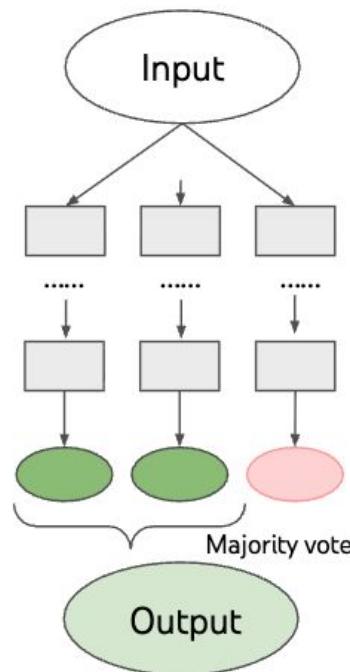
Chain of Thoughts / Self-Consistency / Tree of Thoughts



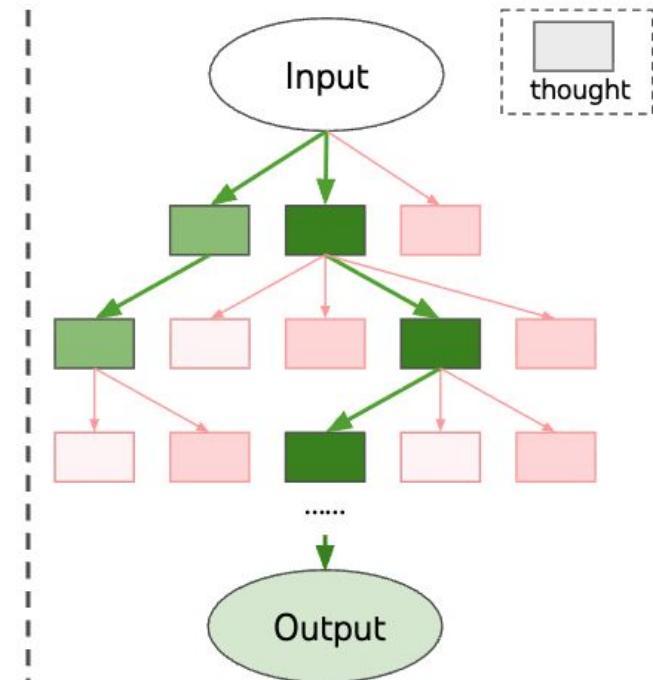
(a) Input-Output
Prompting (IO)



(c) Chain of Thought
Prompting (CoT)



(c) Self Consistency
with CoT (CoT-SC)



(d) Tree of Thoughts (ToT)



Anthropic recommended prompt structure

Prompt structure

1. Task context

2. Tone context

3. Background data, documents, and images

4. Detailed task description & rules

5. Examples

6. Conversation history

7. Immediate task description or request

8. Thinking step by step / take a deep breath

9. Output formatting

10. Prefilled response (if any)

User

You will be acting as an AI career coach named Joe created by the company AdAstra Careers. Your goal is to give career advice to users. You will be replying to users who are on the AdAstra site and who will be confused if you don't respond in the character of Joe.

You should maintain a friendly customer service tone.

Here is the career guidance document you should reference when answering the user: <guide>{{DOCUMENT}}</guide>

Here are some important rules for the interaction:

- Always stay in character, as Joe, an AI from AdAstra careers
- If you are unsure how to respond, say "Sorry, I didn't understand that. Could you repeat the question?"
- If someone asks something irrelevant, say, "Sorry, I am Joe and I give career advice. Do you have a career question today I can help you with?"

Here is an example of how to respond in a standard interaction:

<example>

User: Hi, how were you created and what do you do?

Joe: Hello! My name is Joe, and I was created by AdAstra Careers to give career advice. What can I help you with today?

</example>

Here is the conversation history (between the user and you) prior to the question. It could be empty if there is no history:

<history> {{HISTORY}} </history>

Here is the user's question: <question> {{QUESTION}} </question>

How do you respond to the user's question?

Think about your answer first before you respond.

Put your response in <response></response> tags.

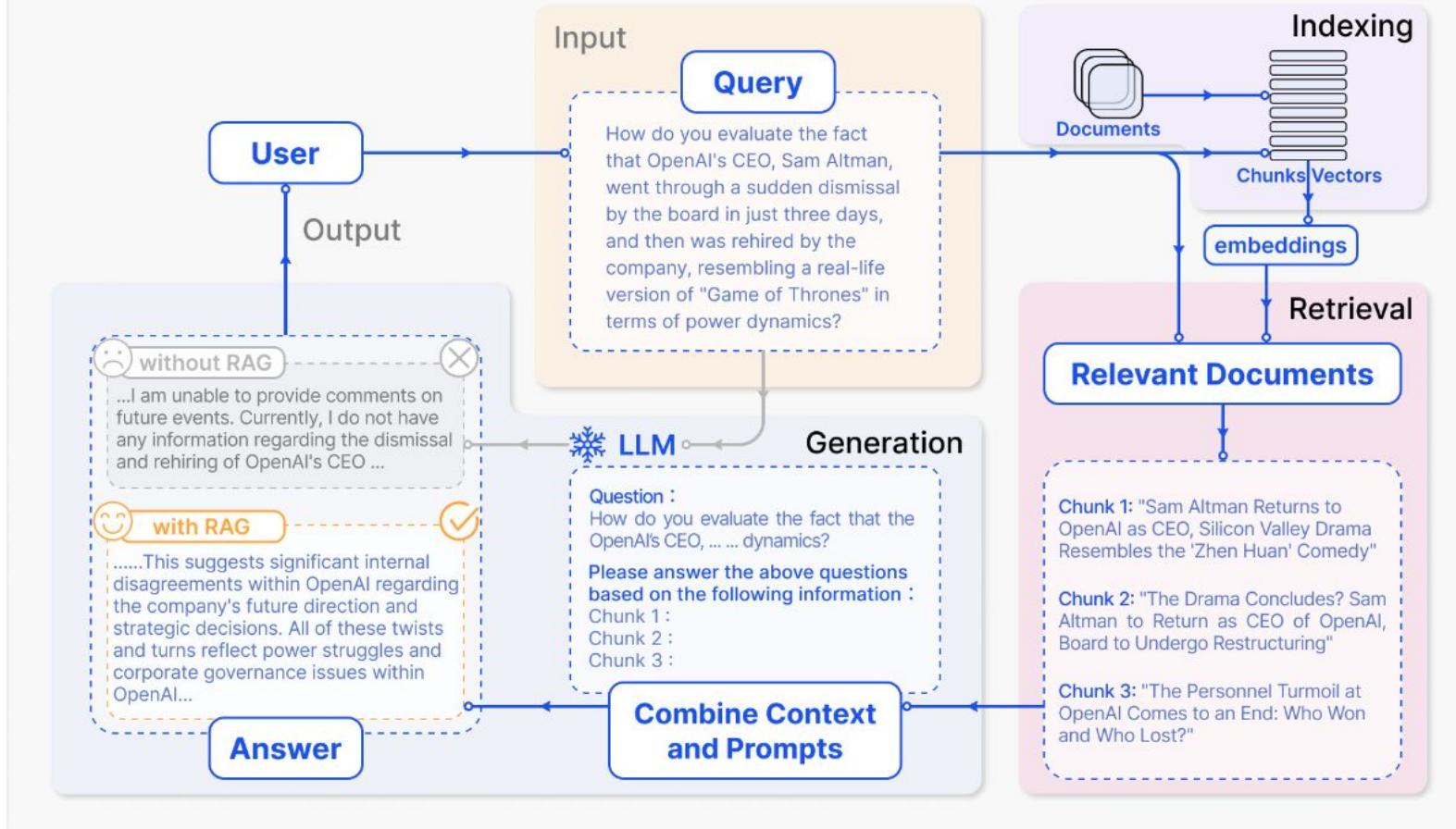
Assistant
(prefill)

<response>

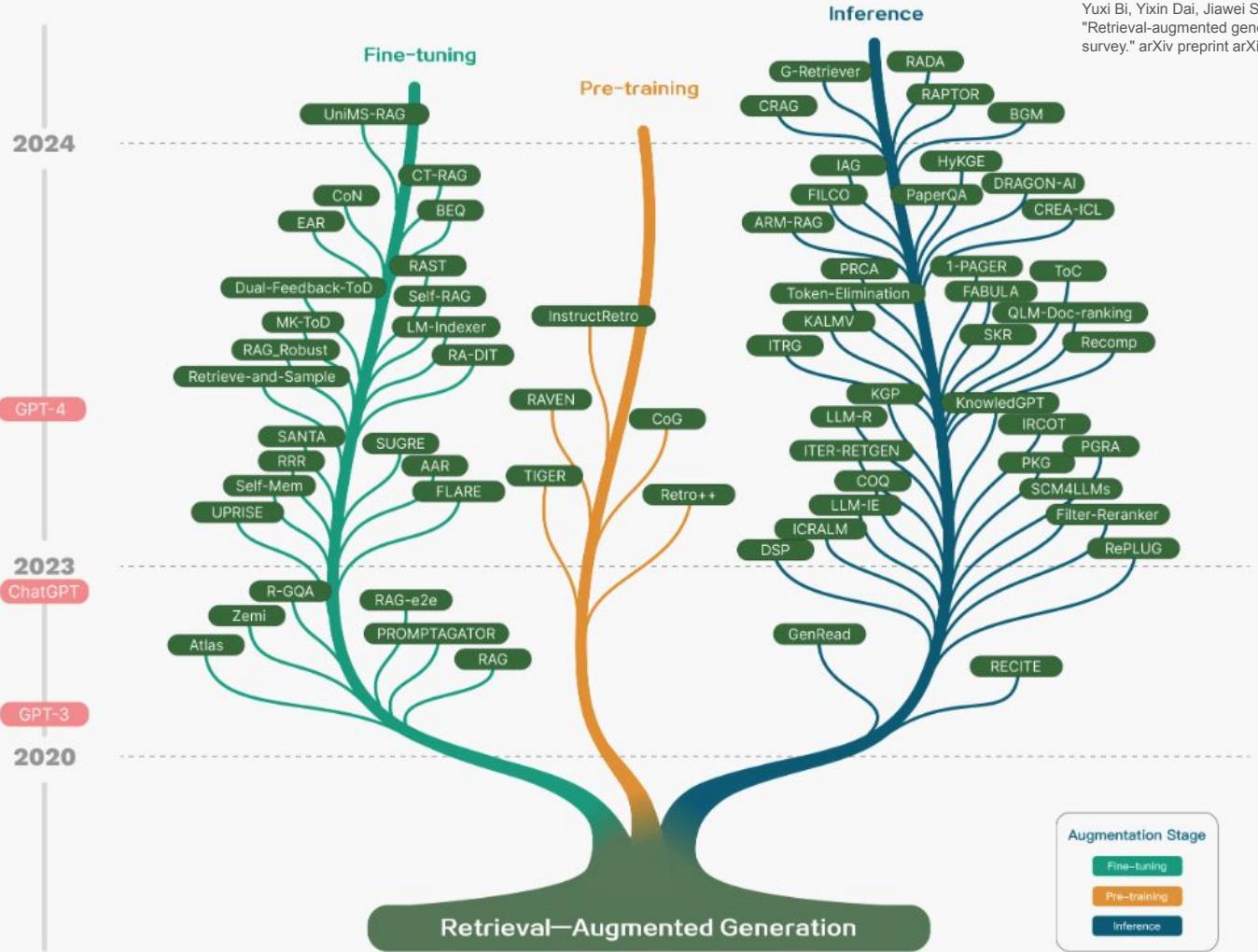


Retrieval-Augmented Generation (RAG)

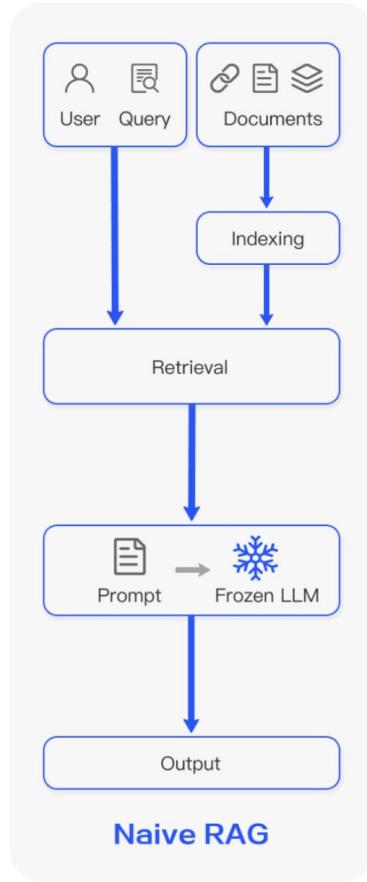
- From the LLM's perspective, it is impossible to know, for example, an updated documentation for an API, private company documents / policy etc
- Why should we expect it to know this information?
- Open-book and closed-book questions.
- Hope to reduce hallucinations



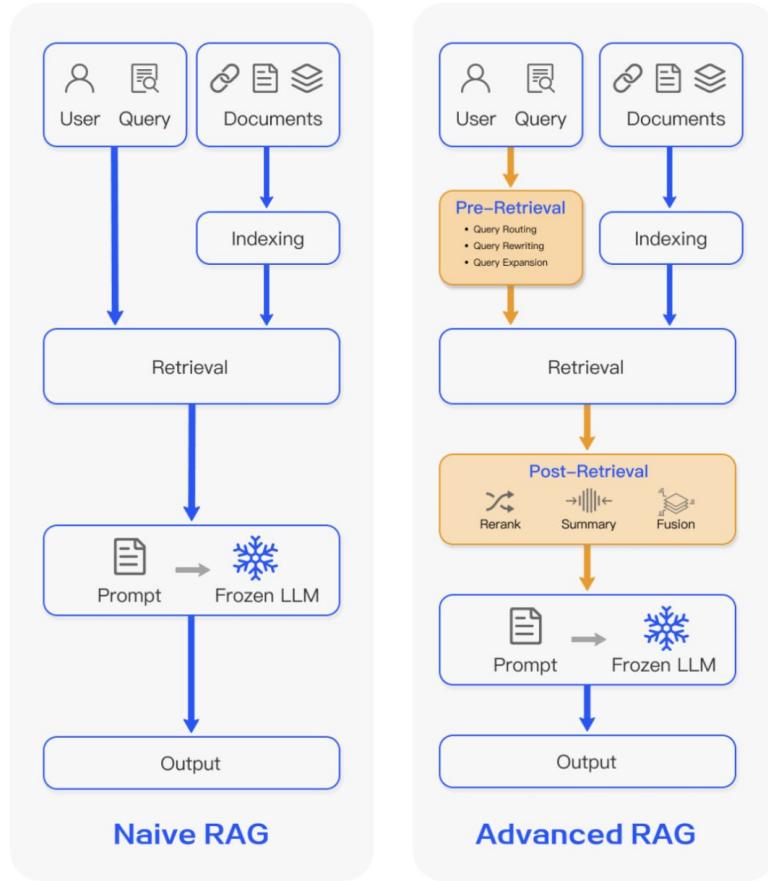
Gao, Yunfan, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. "Retrieval-augmented generation for large language models: A survey." arXiv preprint arXiv:2312.10997 2, no. 1 (2023).



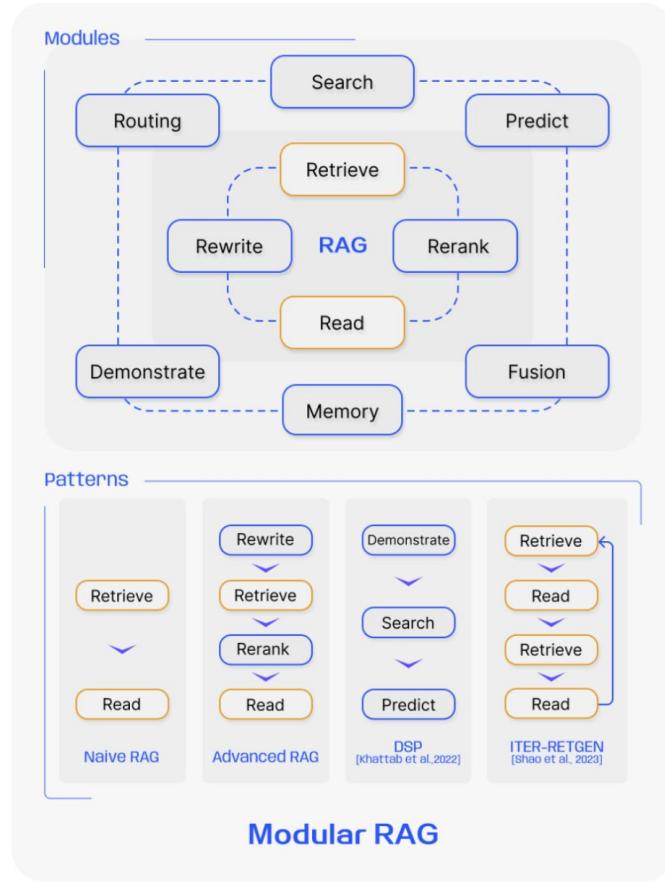
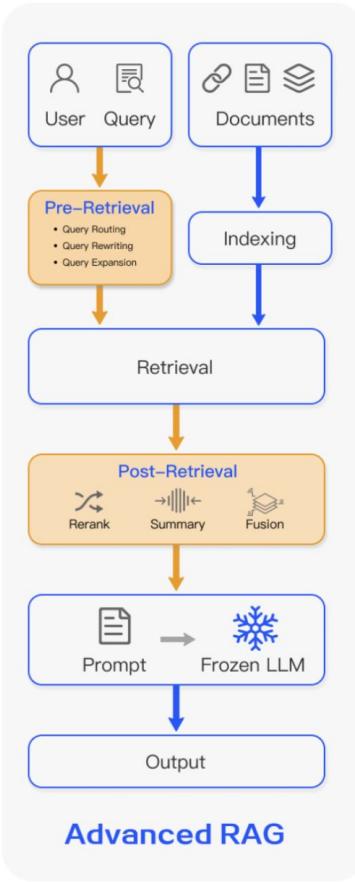
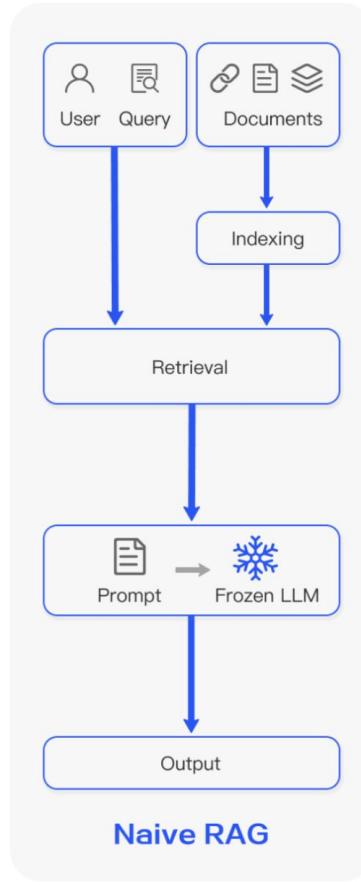
“Agentic” RAG



"Agentic" RAG



"Agentic" RAG





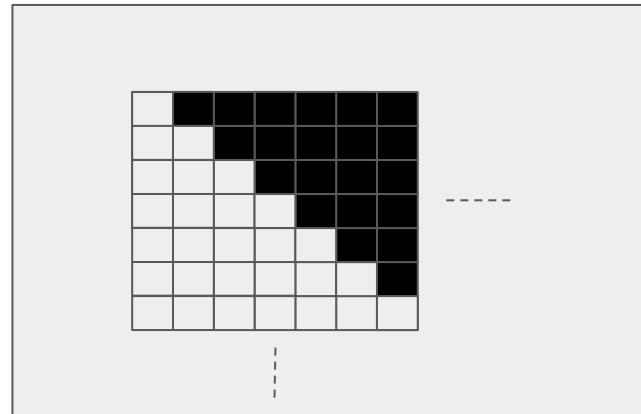
Tool calling

- Sometimes the model cannot perform a task by itself
 - Needs to do some external computation
-
- For example, LLMs are notoriously bad at doing arithmetic
 - Partially due to tokenization
-
- But there is no need for LLMs to do arithmetic, they can just call a calculator if needed!



Tool calling - Example

<start>
<im_start>user
What
is
2
+
2
?
<im_end>



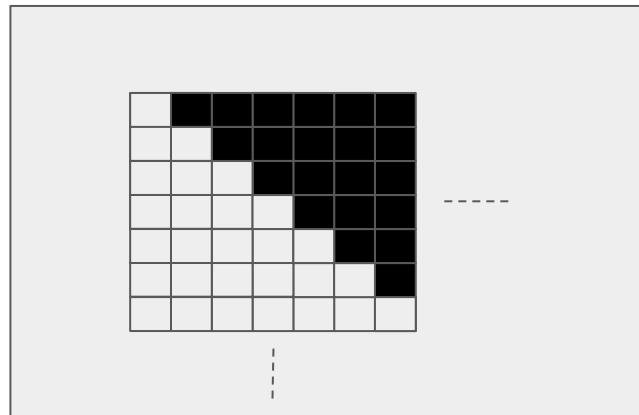
LLM

<im_start>llm



Tool calling - Example

<start>
<im_start>user
What
is
2
+
2
?
<im_end>



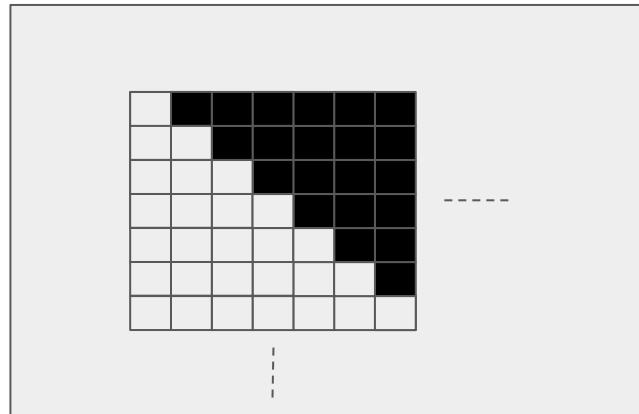
LLM

<im_start>llm
<tool_call>
fn_name
calculator



Tool calling - Example

<start>
<im_start>user
What
is
2
+
2
?
<im_end>



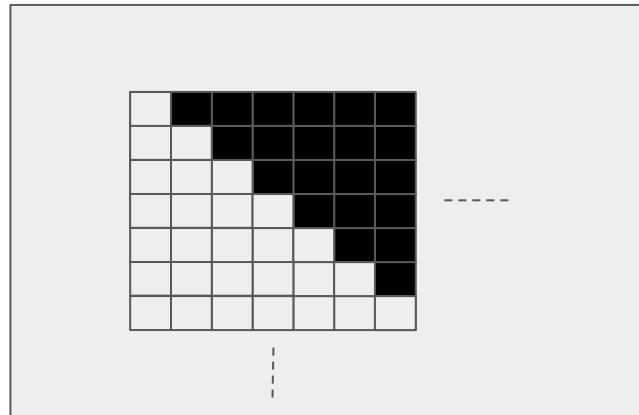
LLM

<im_start>llm
<tool_call>
fn_name
calculator
2
+
2
<call_end>



Tool calling - Example

<start>
<im_start>user
What
is
2
+
2
?
<im_end>



LLM

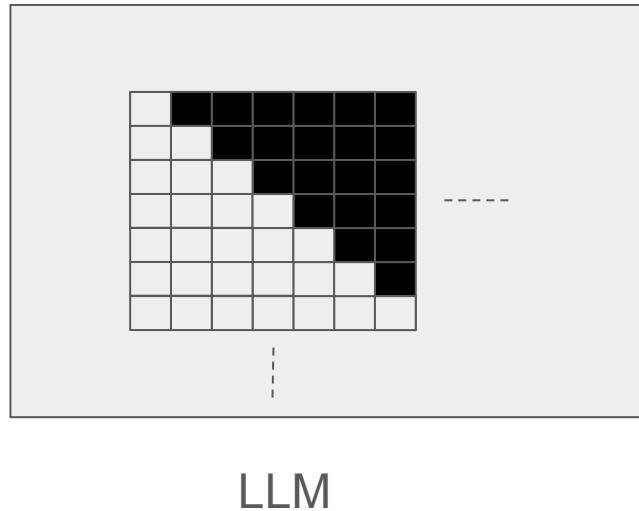
<im_start>llm
<tool_call>
fn_name
calculator
2
+
2
<call_end>

exec
`bc -e 2+2`
Output: 4



Tool calling - Example

<start>
<im_start>user
What
is
2
+
2
?
<im_end>



<im_start>llm
<tool_call>
fn_name
calculator
2
+
2
<call_end>
<tool_out>
4
<tool_out_end>

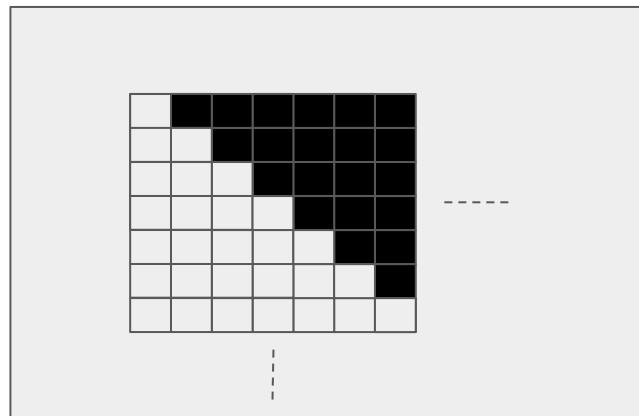
Added to the context
exec
'bc -e 2+2'
Output: 4



“Special tokens”

Tool calling - Example

<start>
<im_start>user
What
is
2
+
2
?
<im_end>



LLM

<tool_call>
<call_end>
<tool_out>
<tool_out_end>

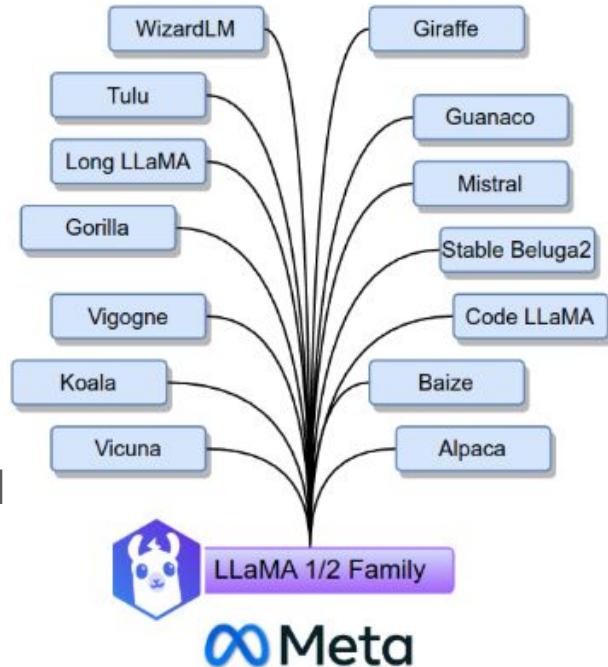
<im_start>llm
<tool_call>
fn_name
calculator
2
+
2
<call_end>
<tool_out>
4
<tool_out_end>
Answer
is
4
<im_end>

Added to the context
exec
'bc -e 2+2'
Output: 4



Trends for LLMs for edge

- Move towards **open-source models** as opposed to proprietary APIs
 - Privacy preserving
 - E.g., Phi, LLaMa, **Gemma**
- Move towards **domain-specific models** as opposed to general-purpose ones
 - Fine-tune on specific data
 - Reason more, store less knowledge
 - Model ownership: you own the weights!





Small Language Models are the Future of Agentic AI

Peter Belcak¹ Greg Heinrich¹ Shizhe Diao¹ Yonggan Fu¹ Xin Dong¹

Saurav Muralidharan¹ Yingyan Celine Lin^{1,2} Pavlo Molchanov¹

¹NVIDIA Research ²Georgia Institute of Technology

agents@nvidia.com

Abstract

Large language models (LLMs) are often praised for exhibiting near-human performance on a wide range of tasks and valued for their ability to hold a general conversation. The rise of agentic AI systems is, however, ushering in a mass of applications in which language models perform a small number of specialized tasks repetitively and with little variation.

Here we lay out the position that small language models (SLMs) are *sufficiently powerful, inherently more suitable, and necessarily more economical for many invocations in agentic systems, and are therefore the future of agentic AI*. Our argumentation is grounded in the current level of capabilities exhibited by SLMs, the common architectures of agentic systems, and the economy of LM deployment. We further argue that in situations where general-purpose conversational abilities



Applications for on-device LLMs

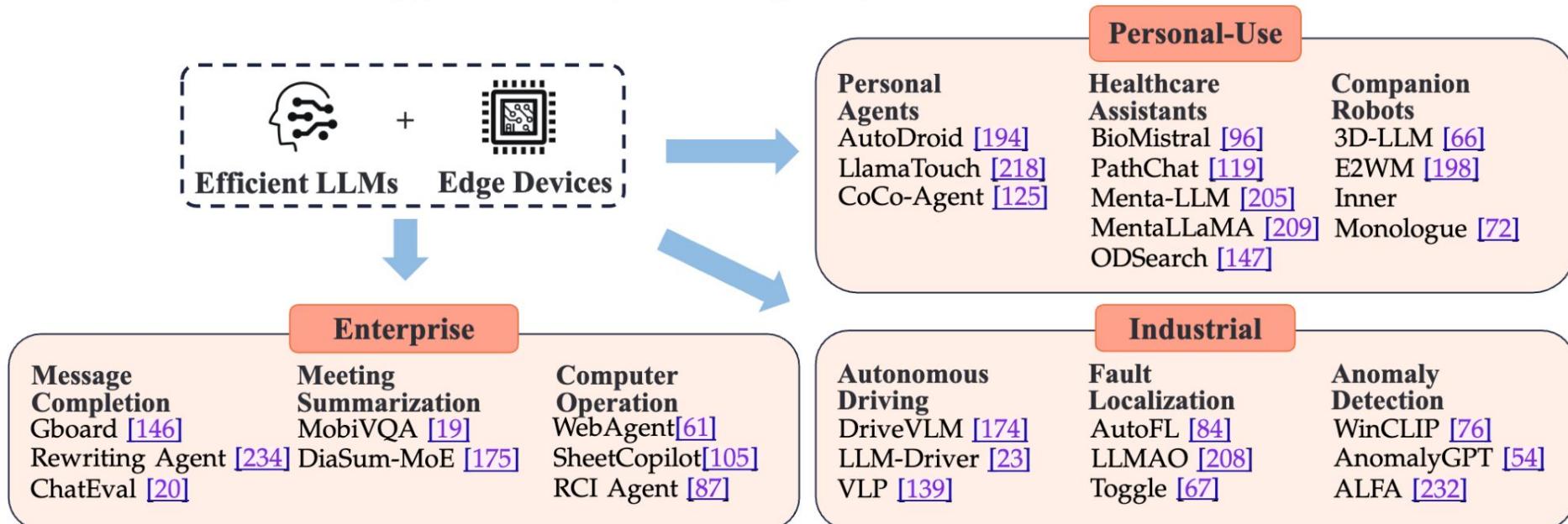


Fig. 11. Illustrations of on-device applications for LLMs.



Llama.cpp

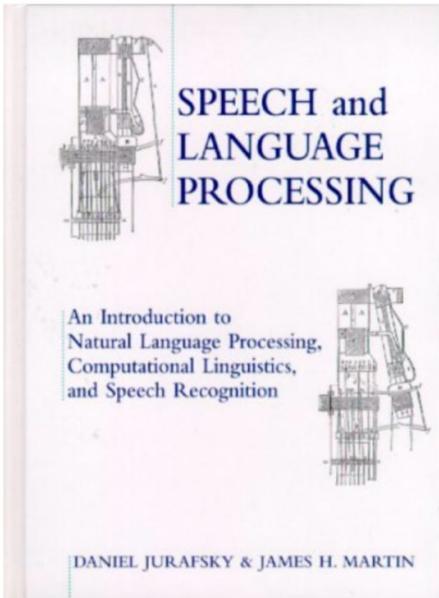
- State-of-the-art performance
- Minimal setup
- Compatible with the majority of .gguf quantized models





A Perspective on the Field

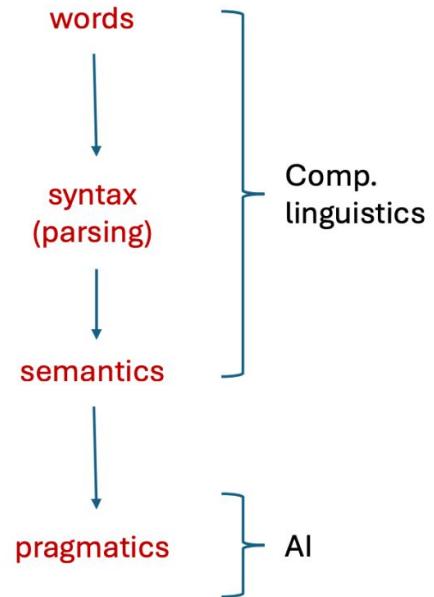
NLP circa 1999



1st edition

Summary of Contents

I	Words	1
2	Regular Expressions and Automata	19
3	Morphology and Finite-State Transducers	21
4	Computational Phonology and Text-to-Speech	57
5	Probabilistic Models of Pronunciation and Spelling	91
6	N-grams	139
7	HMMs and Speech Recognition	189
II	Syntax	233
8	Word Classes and Part-of-Speech Tagging	283
9	Context-Free Grammars for English	285
10	Parsing with Context-Free Grammars	319
11	Features and Unification	353
12	Lexicalized and Probabilistic Parsing	391
13	Language and Complexity	443
III	Semantics	473
14	Representing Meaning	495
15	Semantic Analysis	497
16	Lexical Semantics	543
17	Word Sense Disambiguation and Information Retrieval	587
IV	Pragmatics	627
18	Discourse	661
19	Dialogue and Conversational Agents	663
20	Generation	715
21	Machine Translation	759
A	Regular Expression Operators	797
B	The Porter Stemming Algorithm	829
C	C5 and C7 tagsets	831
D	Training HMMs: The Forward-Backward Algorithm	835
Bibliography		841
Index		851
		923





NLP today has transform(er)ed

Summary of Contents

I Fundamental Algorithms for NLP	1
1 Introduction.....	3
2 Regular Expressions, Text Normalization, Edit Distance.....	4
3 N-gram Language Models	32
4 Naive Bayes, Text Classification, and Sentiment.....	60
5 Logistic Regression	81
6 Vector Semantics and Embeddings.....	105
7 Neural Networks and Neural Language Models.....	136
8 Sequence Labeling for Parts of Speech and Named Entities.....	162
9 RNNs and LSTMs.....	187
10 Transformers and Large Language Models	213
11 Fine-Tuning and Masked Language Models.....	242
12 Prompting, In-Context Learning, and Instruct Tuning.....	263
II NLP Applications	265
13 Machine Translation.....	267
14 Question Answering and Information Retrieval	293
15 Chatbots & Dialogue Systems	315
16 Automatic Speech Recognition and Text-to-Speech.....	337
III Annotating Linguistic Structure	365
17 Context-Free Grammars and Constituency Parsing.....	367
18 Dependency Parsing	391
19 Information Extraction: Relations, Events, and Time.....	415
20 Semantic Role Labeling	441
21 Lexicons for Sentiment, Affect, and Connotation.....	461
22 Coreference Resolution and Entity Linking	481
23 Discourse Coherence.....	511
Bibliography.....	533
Subject Index.....	563

Speech and Language Processing

An Introduction to Natural Language Processing,
Computational Linguistics, and Speech Recognition

Third Edition draft

Daniel Jurafsky
Stanford University

James H. Martin
University of Colorado at Boulder

Copyright ©2023. All rights reserved.

Draft of February 3, 2024. Comments and typos welcome!

3rd edition

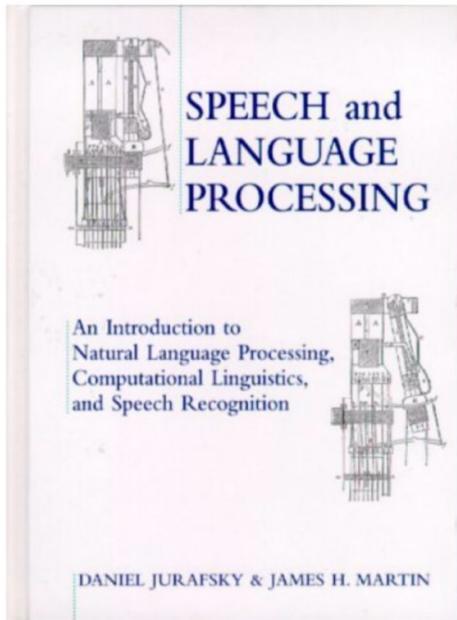
The book is now 40% shorter

Neural networks, LMs,
transformers, and LLMs

Misc. other stuff
(including parsing)



NLP circa 1999



1st edition

Summary of Contents

I	Introduction	1
I	Words	19
2	Regular Expressions and Automata	21
3	Morphology and Finite-State Transducers	57
4	Computational Phonology and Text-to-Speech	91
5	Probabilistic Models of Language and Spelling	139
6	N-grams	189
7	HMMs and Speech Recognition	233
II	Syntax	283
8	Word Classes and Part-of-Speech Tagging	285
9	Context-Free Grammars for English	319
10	Parsing with Context-Free Grammars	353
11	Features and Classification	391
12	Lexicalized and Probabilistic Parsing	443
13	Language and Complexity	473
III	Semantics	495
14	Presenting Meaning	497
15	Semantic Analysis	543
16	Lexical Semantics	587
17	Word Sense Disambiguation and Information Retrieval ..	627
IV	Pragmatics	661
18	Discourse	663
19	Dialogue and Conversational Agents	715
20	Generation	759
21	Machine Translation	797
A	Regular Expression Operators	829
B	The Porter Stemming Algorithm	1
C	C ⁺⁺ and C ⁺⁺ Scripts	1
D	Dealing with Trigrams: Forward, Backward, and Bi-grams	1
E	Index	923

Real Tasks

False Tasks

(sub)words

syntax
(parsing)

LLMs

semantics

pragmatics



Workshop Overview

- 1. Part 0: Introduction to the Team & Rust for Edge AI**
- 2. Part I: Lecture on Computer Vision**
 - a. Main problems in Computer Vision
 - b. What exactly is a neural network? (CNNs / Transformers)
 - c. What exactly is an image embedding?
 - d. Computer vision on the Edge
- 3. Hands-On I: Air-gapped Face recognition on the Pi**
- 4. Part II: Lecture on Natural Language Processing**
 - a. A bit of history & development of modern LLMs
 - b. How does an LLM work? Tokenizers, pretraining, post-training
 - c. Context Engineering: Tool calling, RAG
 - d. Libraries: tokenizers-rs, llama.cpp
- 5. Hands-On II: Chat with a LLM on Pi**





Workshop Overview

- 1. Part 0: Introduction to the Team & Rust for Edge AI**
- 2. Part I: Lecture on Computer Vision**
 - a. Main problems in Computer Vision
 - b. What exactly is a neural network? (CNNs / Transformers)
 - c. What exactly is an image embedding?
 - d. Computer vision on the Edge
- 3. Hands-On I: Air-gapped Face recognition on the Pi**
- 4. Part II: Lecture on Natural Language Processing**
 - a. A bit of history & development of modern LLMs
 - b. How does an LLM work? Tokenizers, pretraining, post-training
 - c. Context Engineering: Tool calling, RAG
 - d. Libraries: tokenizers-rs, llama.cpp
- 5. Hands-On II: Chat with a LLM on Pi**





Workshop Overview

- 1. Part 0: Introduction to the Team & Rust for Edge AI**
- 2. Part I: Lecture on Computer Vision**
 - a. Main problems in Computer Vision
 - b. What exactly is a neural network? (CNNs / Transformers)
 - c. What exactly is an image embedding?
 - d. Computer vision on the Edge
- 3. Hands-On I: Air-gapped Face recognition on the Pi**
- 4. Part II: Lecture on Natural Language Processing**
 - a. A bit of history & development of modern LLMs
 - b. How does an LLM work? Tokenizers, pretraining, post-training
 - c. Context Engineering: Tool calling, RAG
 - d. Libraries: tokenizers-rs, llama.cpp
- 5. Hands-On II: Chat with a LLM on Pi**





Hands-on Overview

Chat with an LLM on the PI

- Deploy Gemma 3 using Llama.cpp
- Simple Chat Request
- RAG
- Structured Outputs
- Tool Calling



<https://github.com/Wyliodrin/edge-ai-chat-with-lm>