# Week 14

## Machine Learning and Big Data - DATA622

CUNY School of Professional Studies

# Meaning of words

Words may have many meanings (polysemy). The meaning of a word depends on its context.
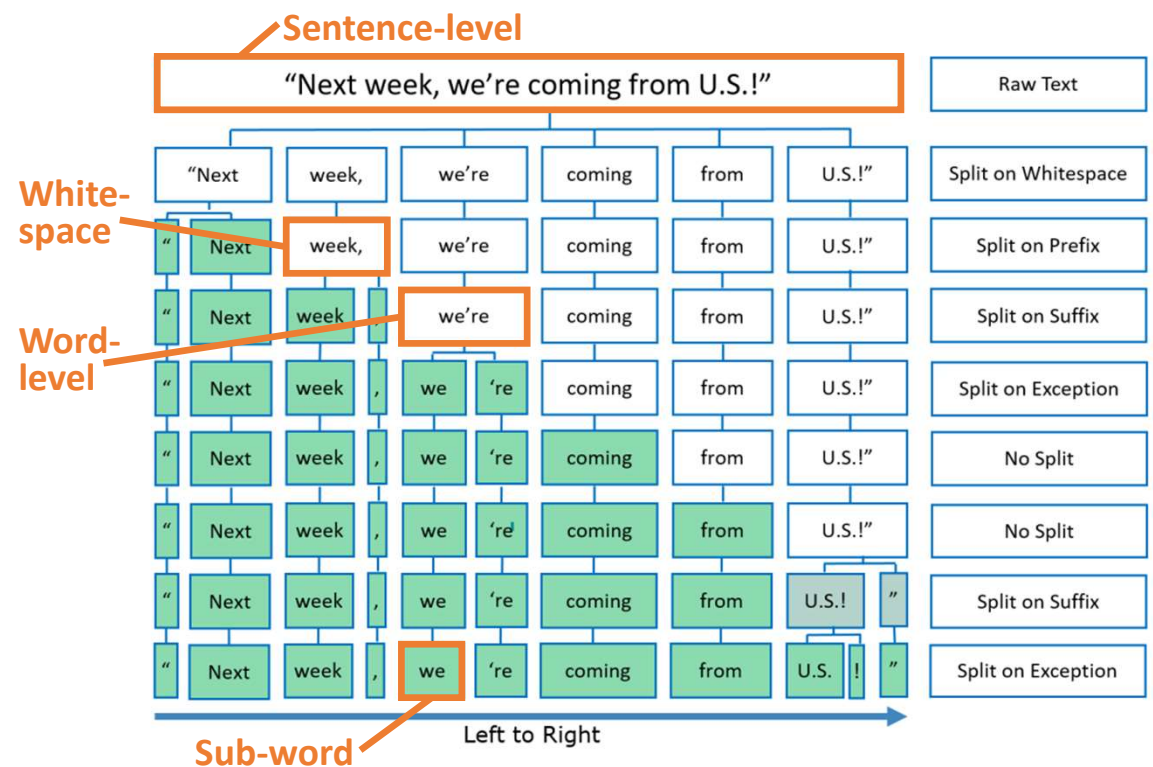
An Example: **bow**
- the front of a ship
- to bend forward in respect
- a weapon that shoots arrows
- to bend outward

CUNY | School of Professional Studies

# Tokenization

Tokenization is the processes of splitting text into manageable pievces: tokens
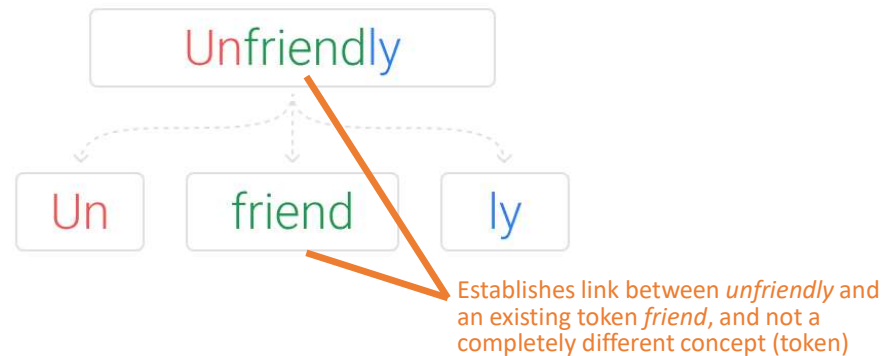
Types of tokenization:

- Character-level
- Subword-level
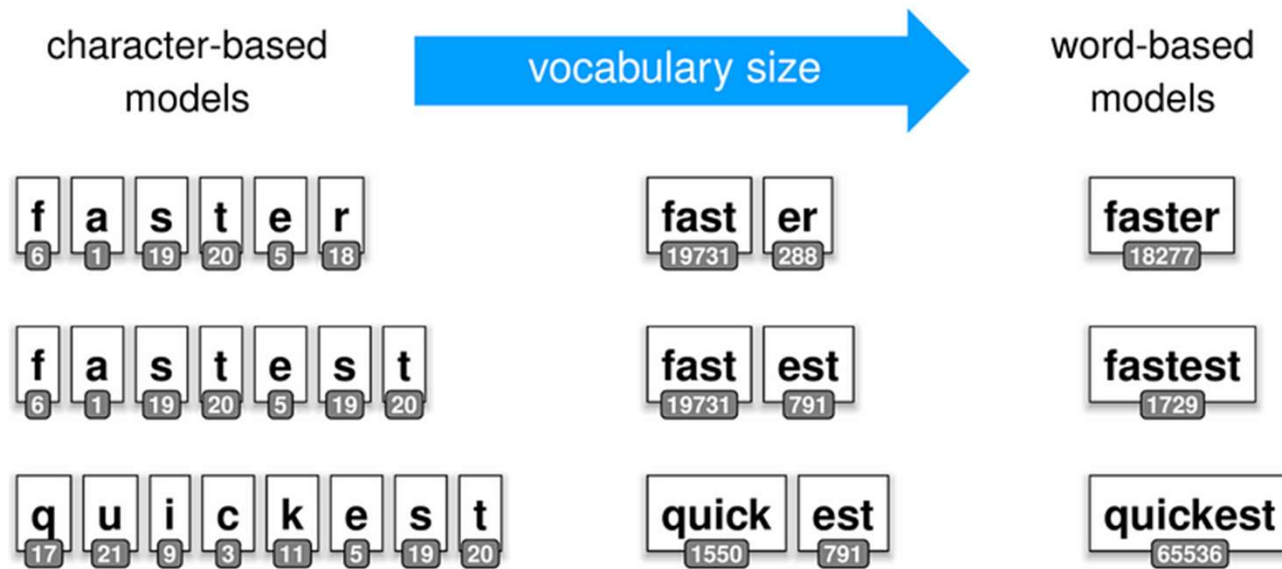- Word-level
- Whitespace-level
- Sentence-level



**Sentence-level**

| | | | | | | | Raw Text |
| --- | --- | --- | --- | --- | --- | --- | --- |

"Next week, we're coming from U.S.!"

**White-space**

**Word-level**

**Sub-word**

Left to Right

# Tokenization

Sub-word is most popular, best balance
of vocabulary size and retention of meaning



Unfriendly

Un    friend    ly

Establishes link between *unfriendly* and
an existing token *friend*, and not a
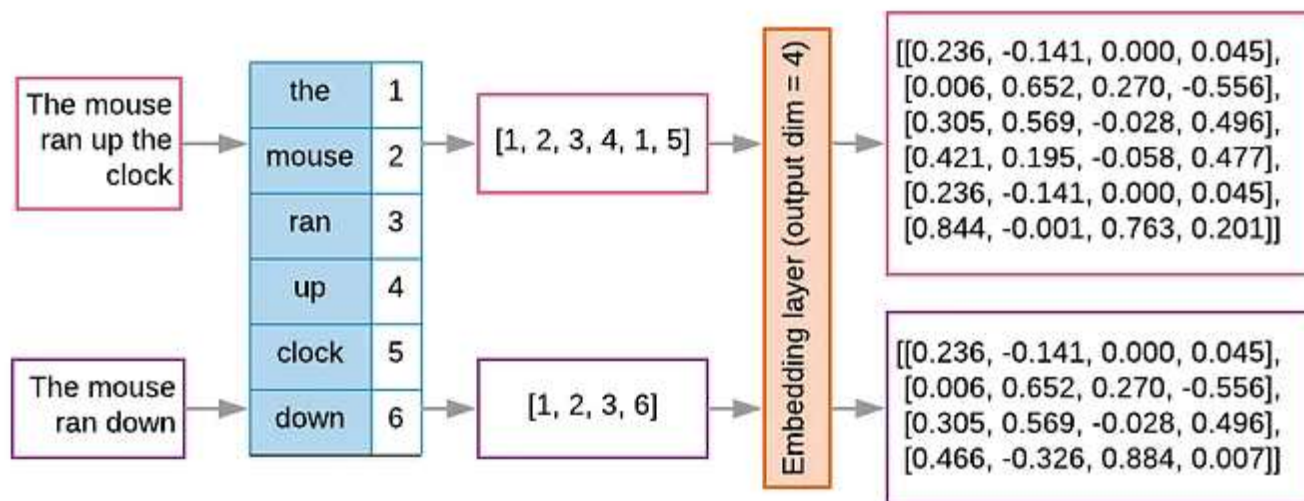completely different concept (token)

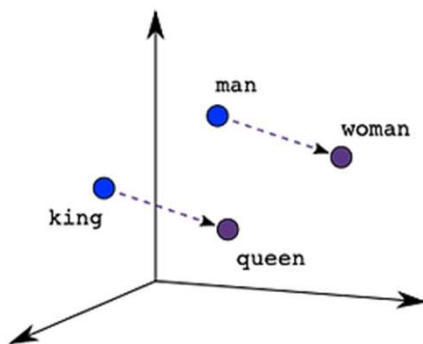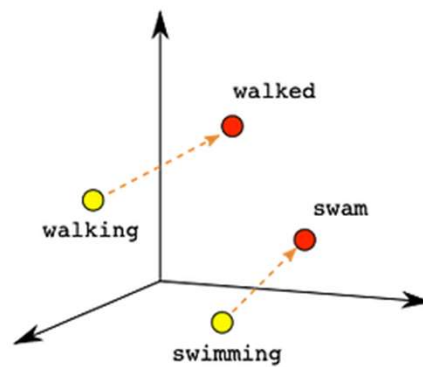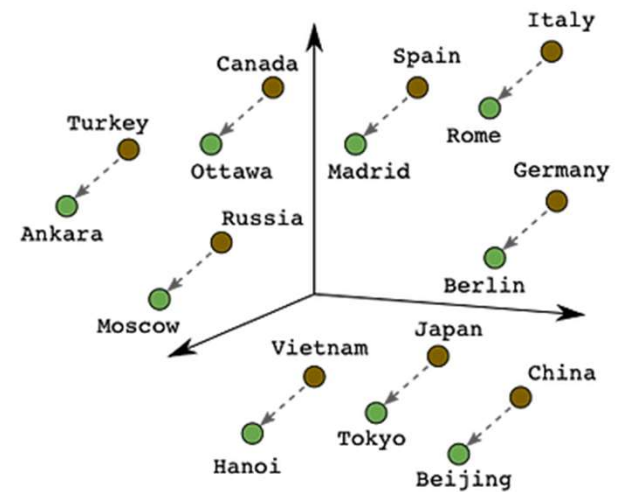CU | School of
NY | Professional Studies

# Tokenization

# Embeddings

# Embeddings



Male-Female
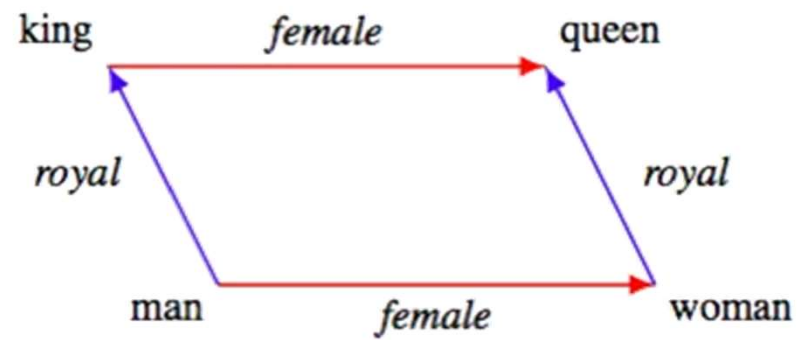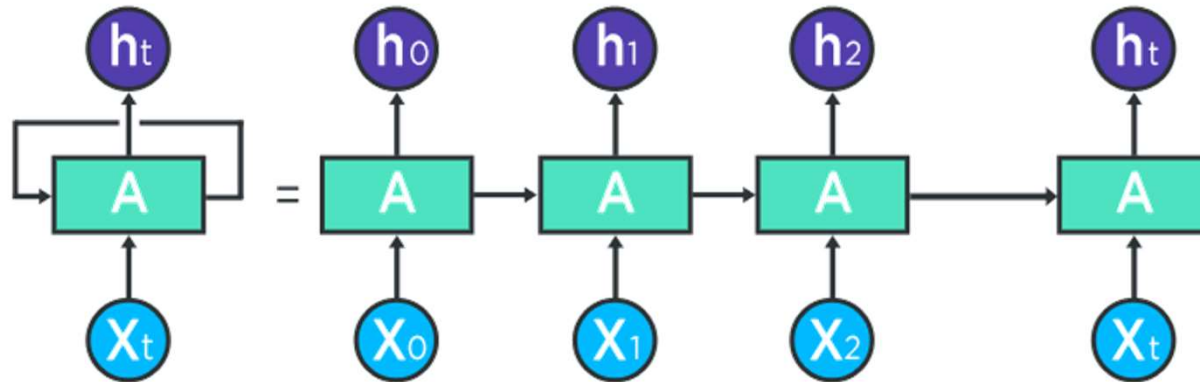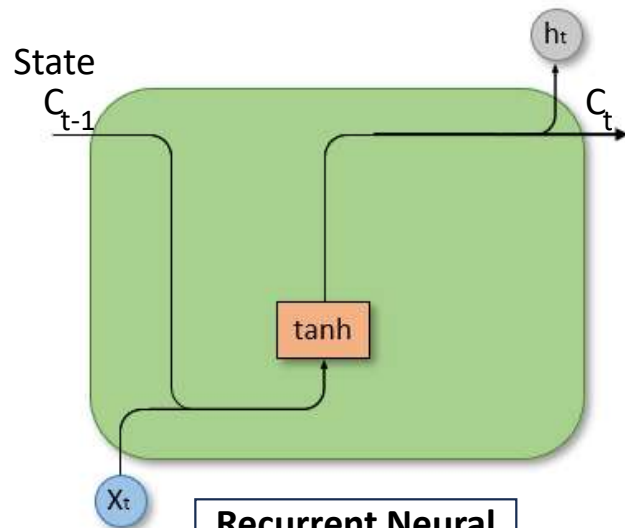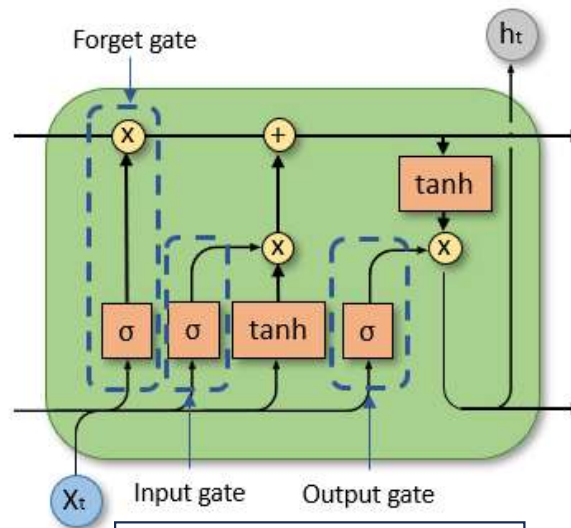
Verb Tense

Country-Capital

# Embeddings
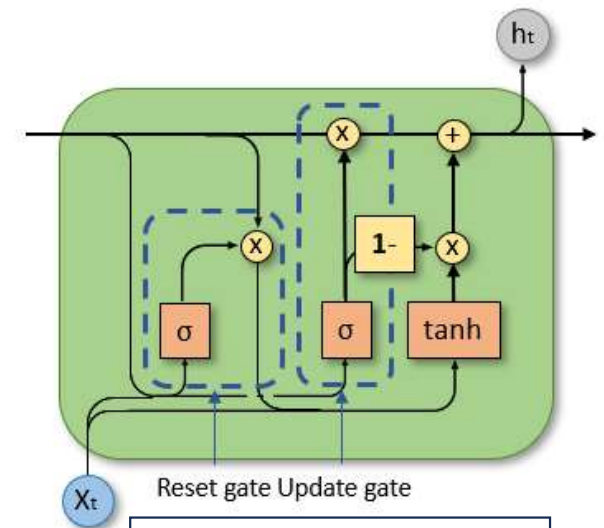
# RNN

# Evolution: RNN – LSTM - GRU



**Recurrent Neural Network (RNN)**

**Long Short-Term Memory (LSTM)**
The LSTM adds three gates (forget gate, input gate, and output gate) to the RNN

**Gated recurrent unit (GRU)**
Variant of the LSTM that synthesizes forget and input gate into a single update gate

# Attention

## 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
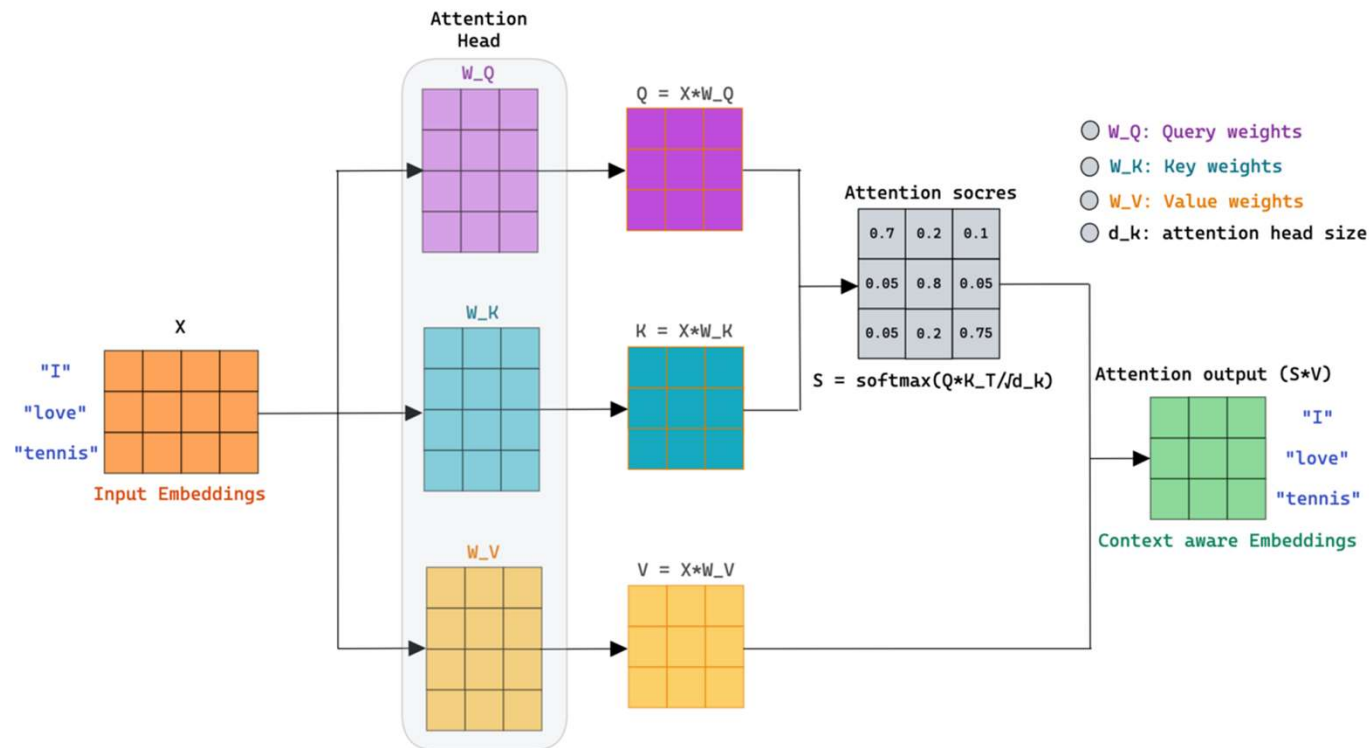
**Jakob Uszkoreit**[*]
Google Research
usz@google.com

**Llion Jones**[*]
Google Research
llion@google.com

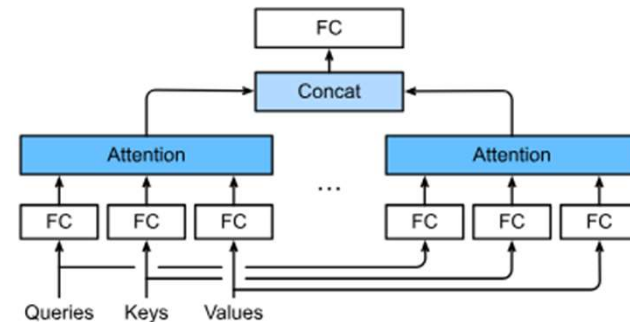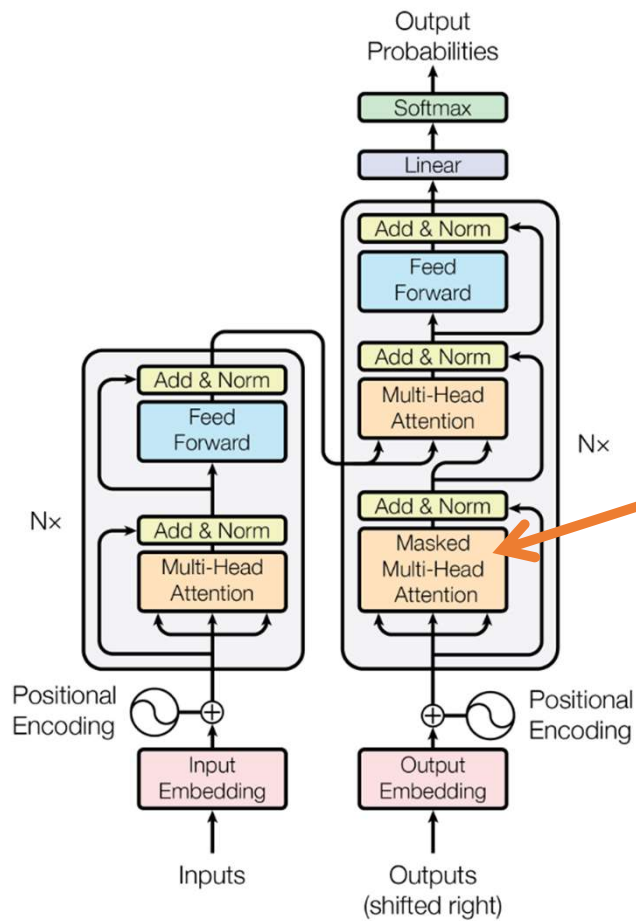**Aidan N. Gomez**[*][†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin**[*][‡]
illia.polosukhin@gmail.com

CUNY | School of Professional Studies

# Attention

# Transformers

# Large Language Models

# Large Language Models

## Large Language Models (LLMs)
- Pre-trained with extremely large datasets – architected to scale
- Can be adopted to a wide range of downstream tasks
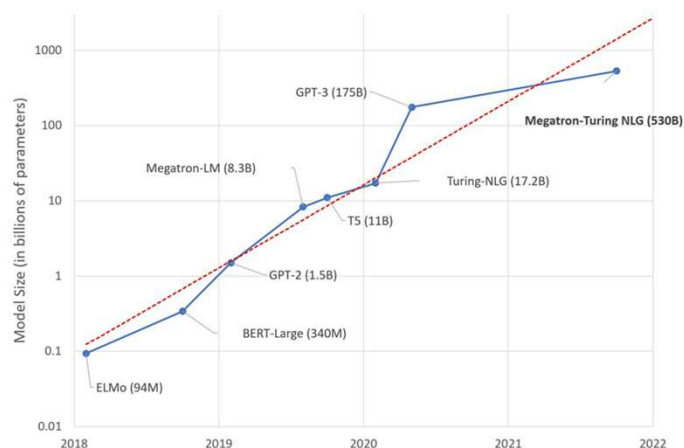- LLMs are Foundation Models.

## Very Large training datasets

**GPT-3 Datasets Summary.**

|  | Wikipedia | Books | Journals | Reddit links | CC | Other | Total |
|---|---|---|---|---|---|---|---|
| GB | 11.4 | 21 | 101 | 50 | 570 | | 753GB |
| Tokens | 3 | 12 | 55 | 19 | 410 | | 499BTokens |

## LLM Scaling: a new Moore's Law



CUNY | School of Professional Studies

# Challenges of ML Today

**Data**
- Task-specific data
- Feature Engineering
- 'Hand-crafted'

**Models**
- Domain-specific
- Limited re-use
- 'Hand-crafted'

**Talent**
- Data Scientists
- ML Researchers
- Specialized skills

Traditional Machine Learning

**Time**
- Long ML lifecycles
- Hard to do quick tests
- Slow retraining cycle

**Costs**
- High retraining costs
- Limited enterprise scaling
- Hard to prototype until later in the project

16

CUNY | School of Professional Studies

# Benefits of LLMs

**Increased Velocity**
- Focus shifts from training models from scratch, to fine-tune models
- Faster time to market

**Increased Opportunity for AI/ML involvement**
- Potential to scale to a wider pool of users to perform AI/ML
- Simple text interface and natural language instruction

**Cost effective**
- Scale to multiple use-cases per LLM
- Faster / lower-cost prototyping

**Tapping into state-of-the-art AI**
- Few-shot learners (and 'Surprisingly good without fine tuning')
- Perform tasks not explicitly trained on

**Emergent Capabilities**
- Emergent capabilities that surface with LLM size
  - i.e. capabilities not present in smaller models but emerge in larger models
- LLM scale highly correlated with downstream performance[1]

17

# Benefits of LLMs



**Fine Tuning**
- Load foundational model
- Add task-specific prompts
- Minimal data, compute, time
- SOTA results

**Embeddings**
- Encode Content into dense vector for downstream use
- Use in downstream models or similarity search

**In-context Learning**
- "Ask" the model to perform a task as part of the input
- Provide examples to help

# Risks of LLMs

**Bias Propagation**
- Potential bias / toxic output
- Responsible AI is critical

**LLMs fail in subtle ways**
- Hallucinations
- Evaluation / safeguards required

**Increasing model scale**
- Exponential growth in size
- Complexity in training / deployment

**LLM costs**
- Cost / Latency trade-off
- Understand & manage costs

**Closed-Source models**
- Legal restrictions to some models



19